

Deep learning for image recognition II

Data Science Retreat Berlin | 28.02.2017
Ludwig Schmidt-Hackenberg

Quiz I

Q: What does the weights vector do?

A: It defines the importance of a neuron's inputs

Q: What does the bias do?

A: It changes the neuron's output independent of the input

Q: What is gradient descent?

A: A optimization algorithm to find a local minimum for a function

Quiz II

Q: What is the learning rate?

A: Is the step size (factor) of the gradient descent

Q: What is stochastic gradient descent?

A: Using only a random (thus stochastic) subset of the data for gradient descent

Q: What is an epoch?

A: A SGD training run using all data

Quiz III

Q: What is the cost function?

A: A function that returns a number representing how well the neural network performed mapping training examples to correct output.

Q: What is a loss function?

A: A different name for the cost function.

Q: What is regularization?

A: (Here) The process of introducing additional information into a learning process to prevent overfitting.

Quiz III

Q: What is the cost function?

A: A function that returns a number representing how well the neural network performed mapping training examples to correct output.

Q: What is a loss function?

A: A different name for the cost function.

Q: What is regularization?

A: (Here) The process of a learning process to pre-

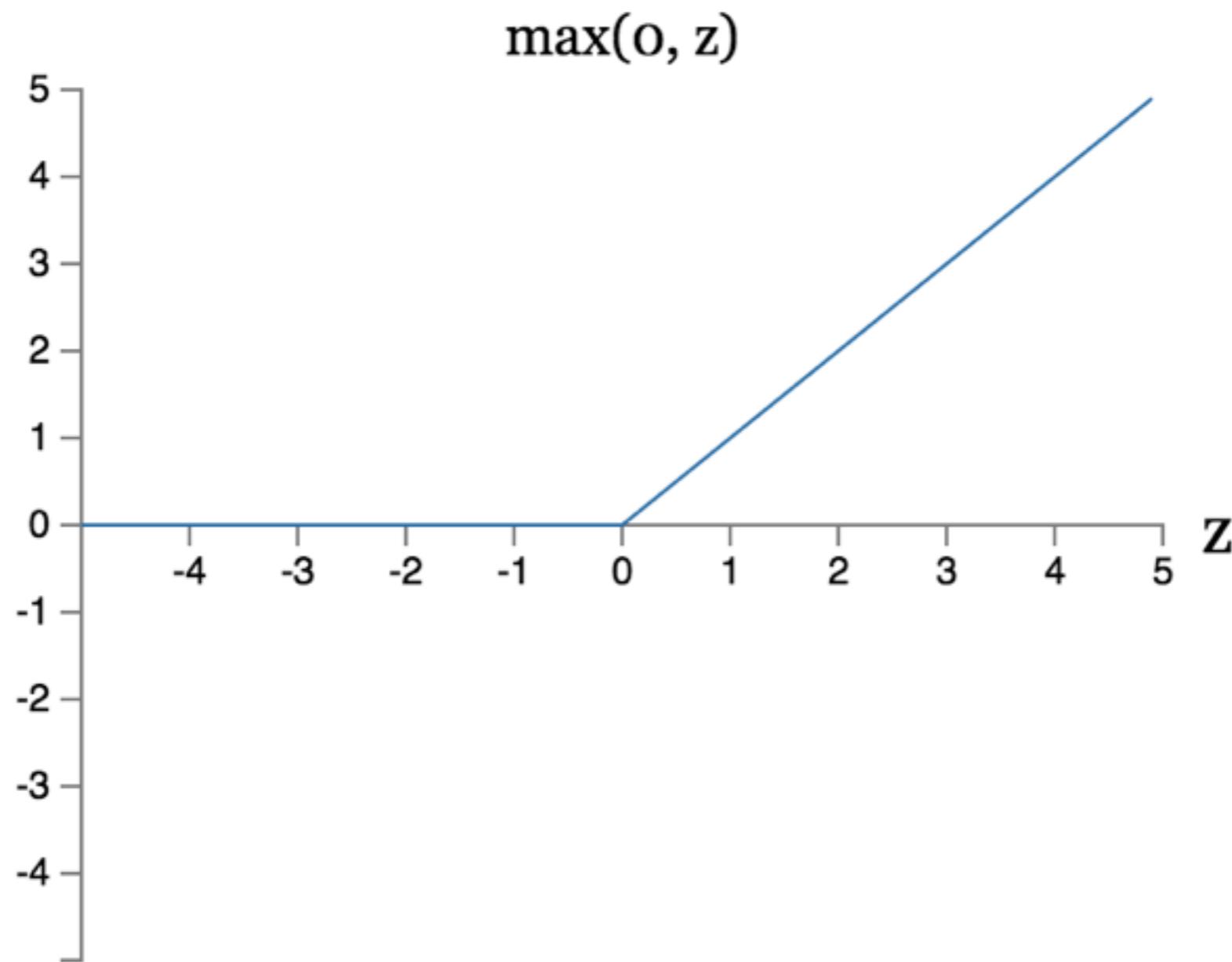
And why now?
GPU and datasets!

Deep Learning II

‘Advanced’ Neural Network techniques II

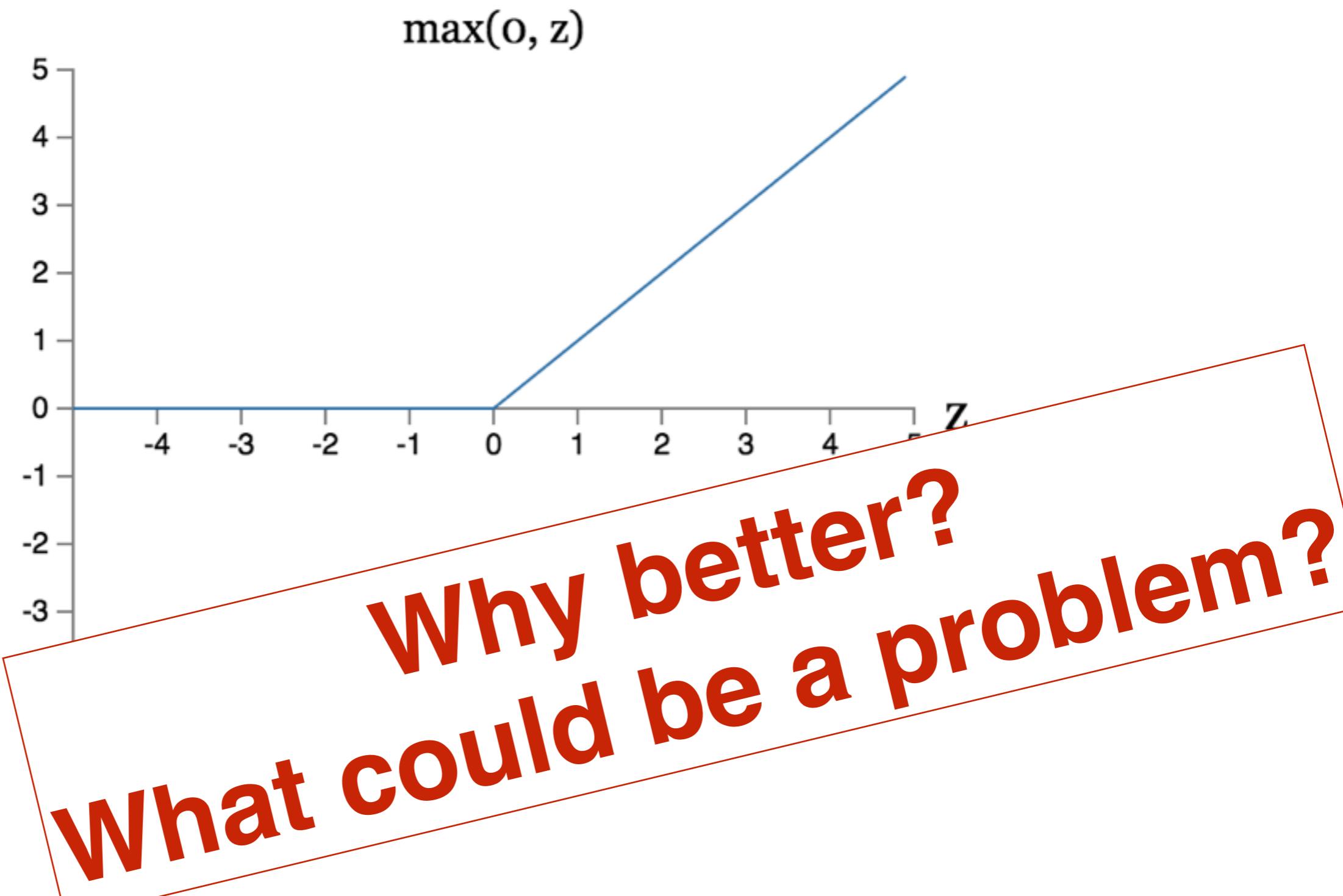
Rectified linear units

A better activation function

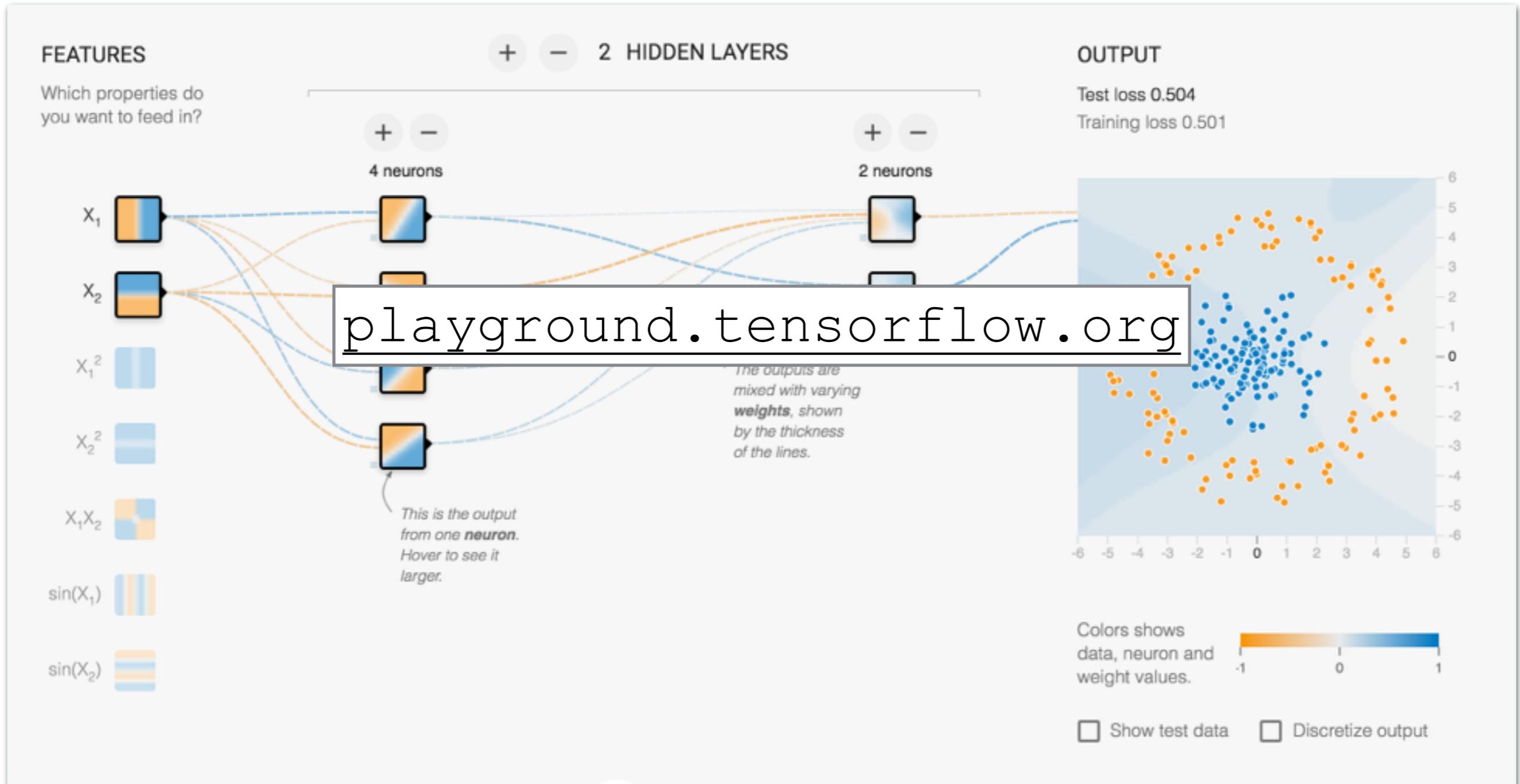


Rectified linear units

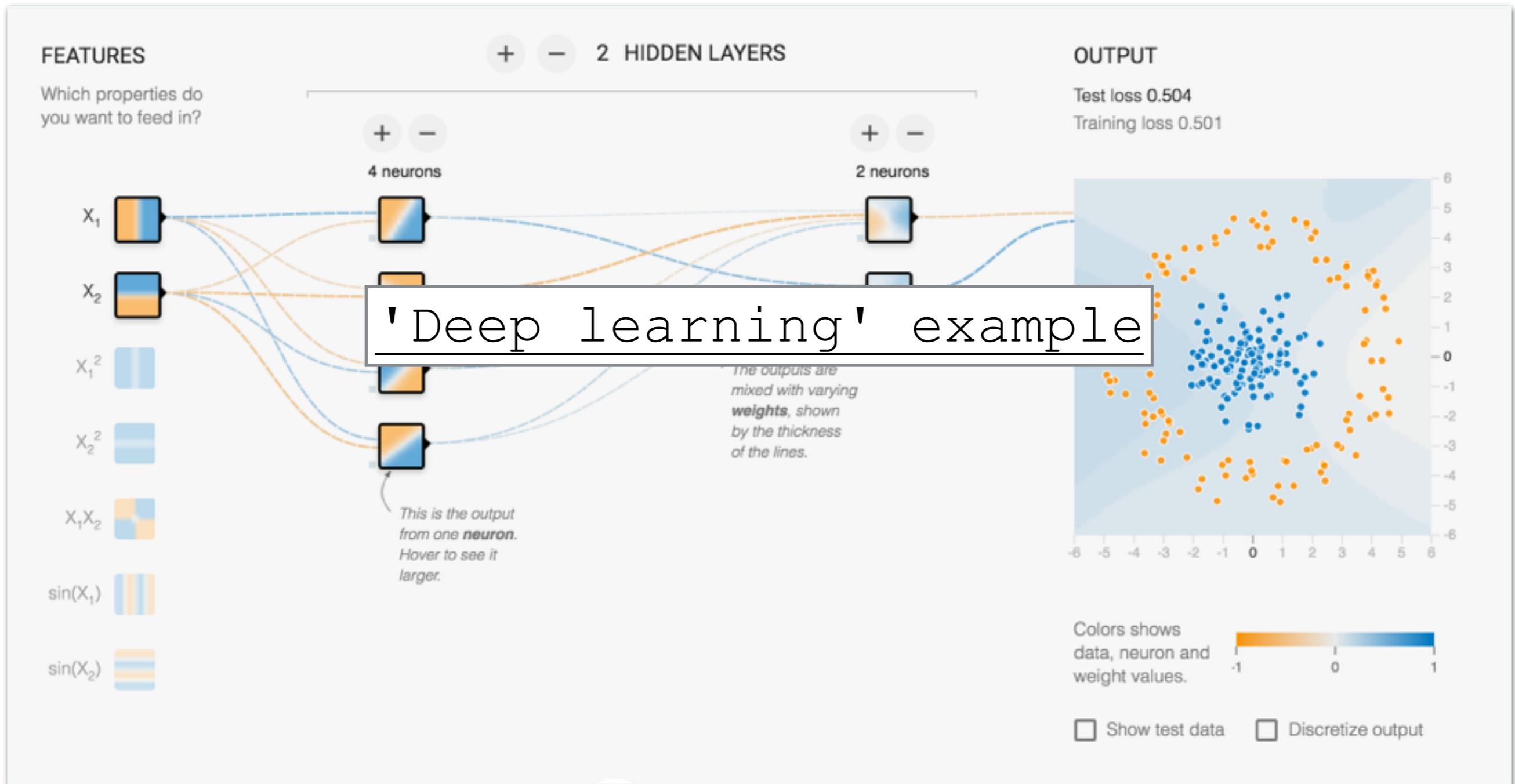
A better activation function



Tensor flow Playground

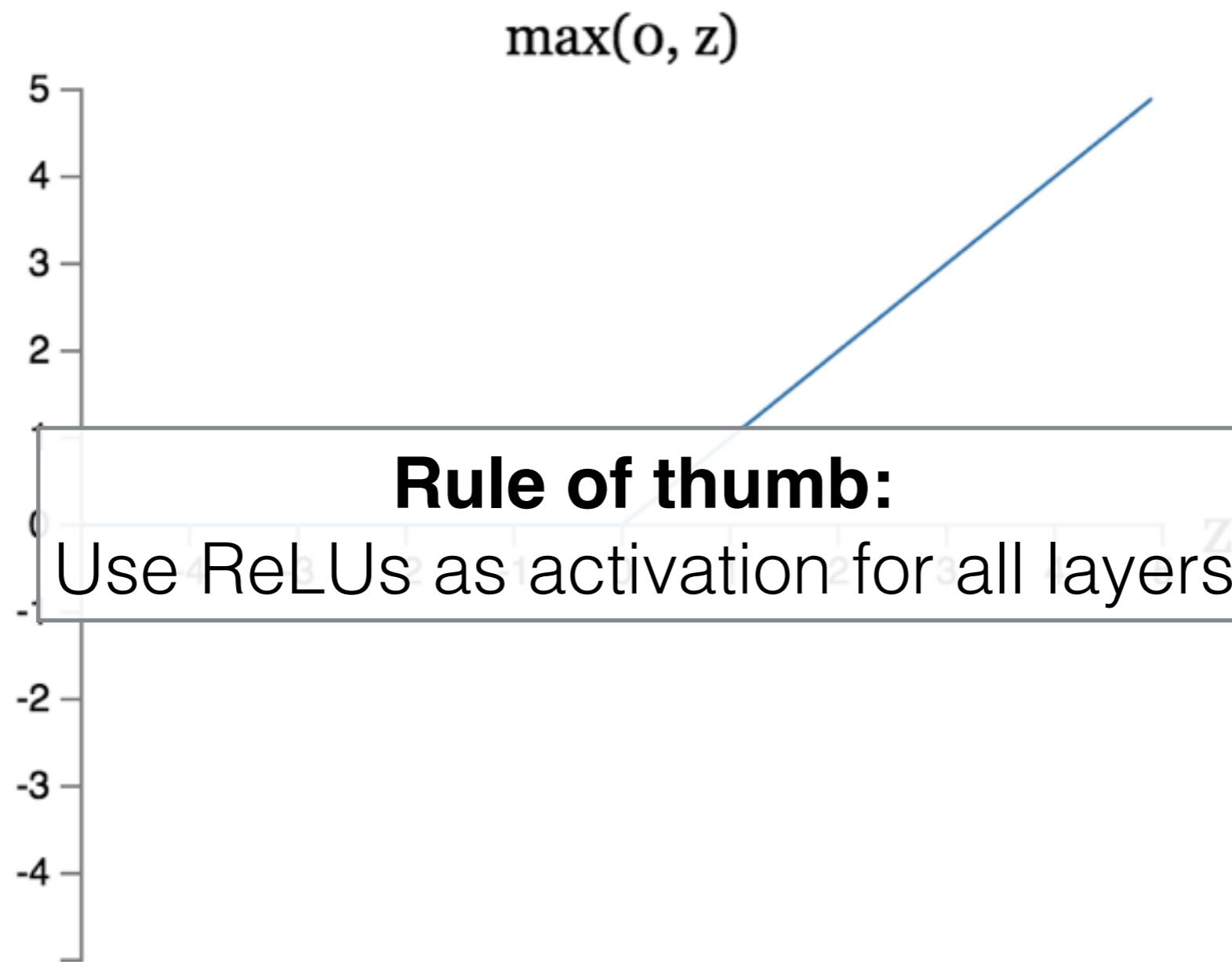


Tensor flow Playground



Rectified linear units

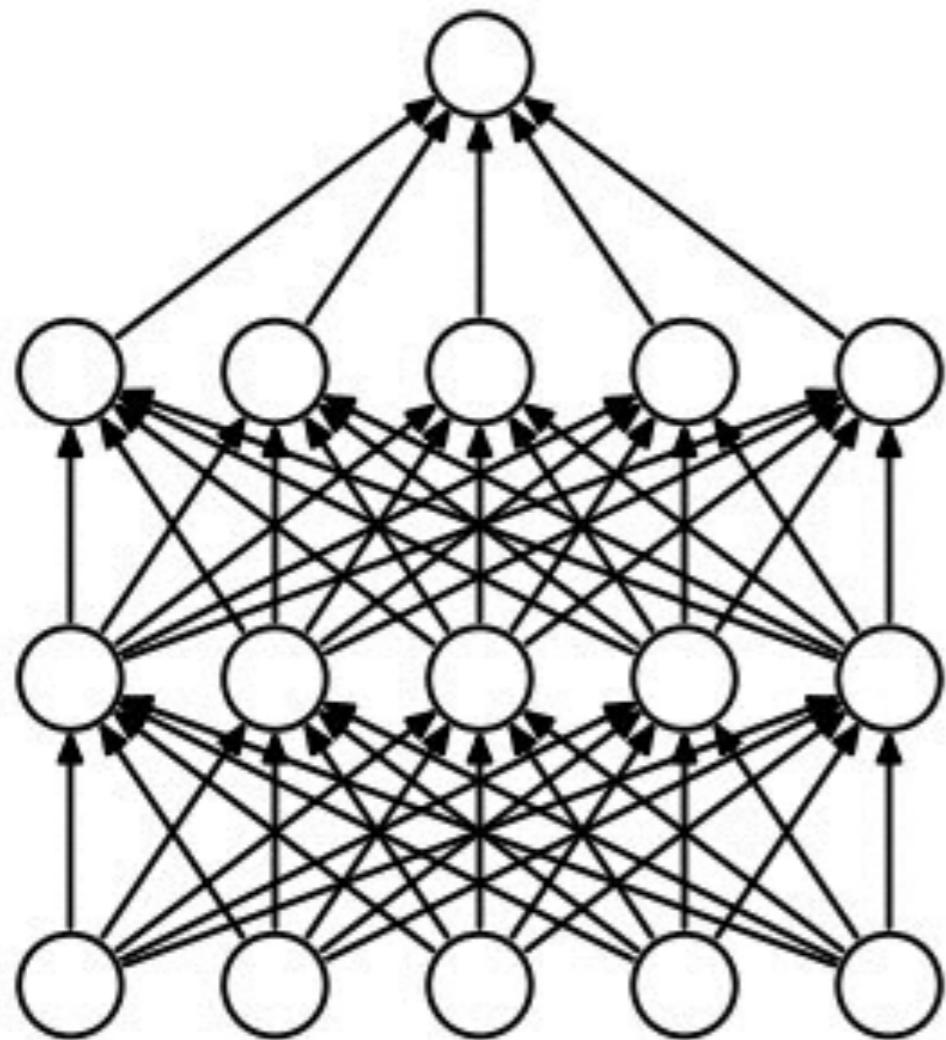
A better activation function



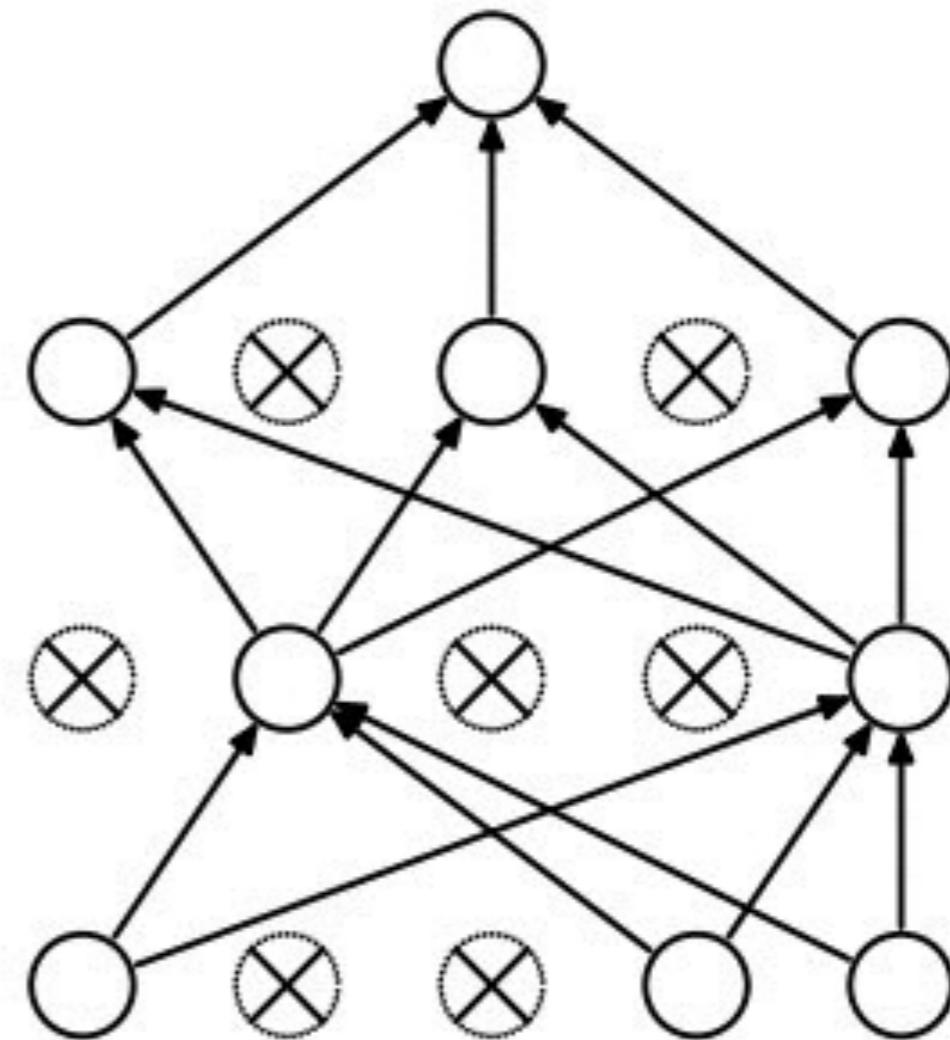
Why does Random
Forests work so well?

Regularization with Dropout

random ensembles for NN



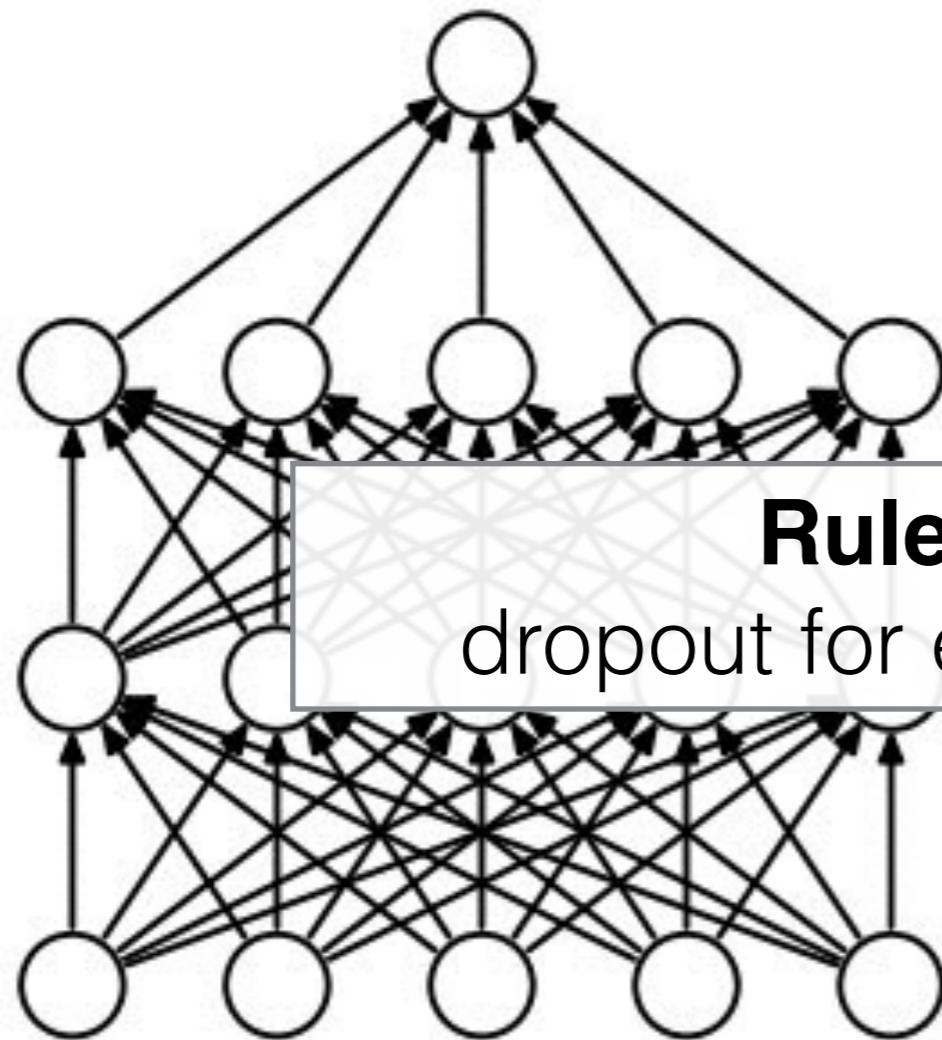
(a) Standard Neural Net



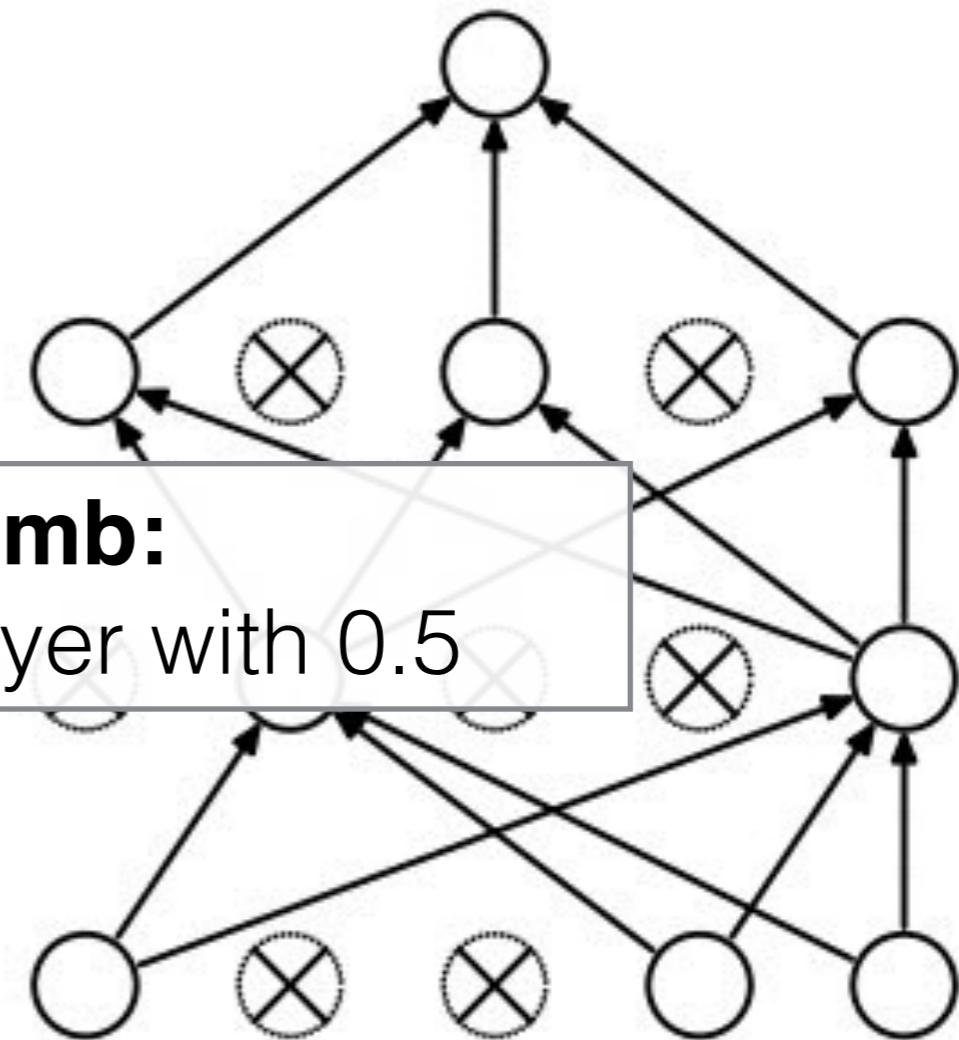
(b) After applying dropout.

Regularization with Dropout

random ensembles for NN



(a) Standard Neural Net



(b) After applying dropout.

Pimp my
Gradient Descent

Decaying the learning rate

- Every n-epochs by a constant factor
- Every step exponentially
- When the learning stalls
-

Gradient Descent

$$w \rightarrow w' = w - \eta \nabla C$$

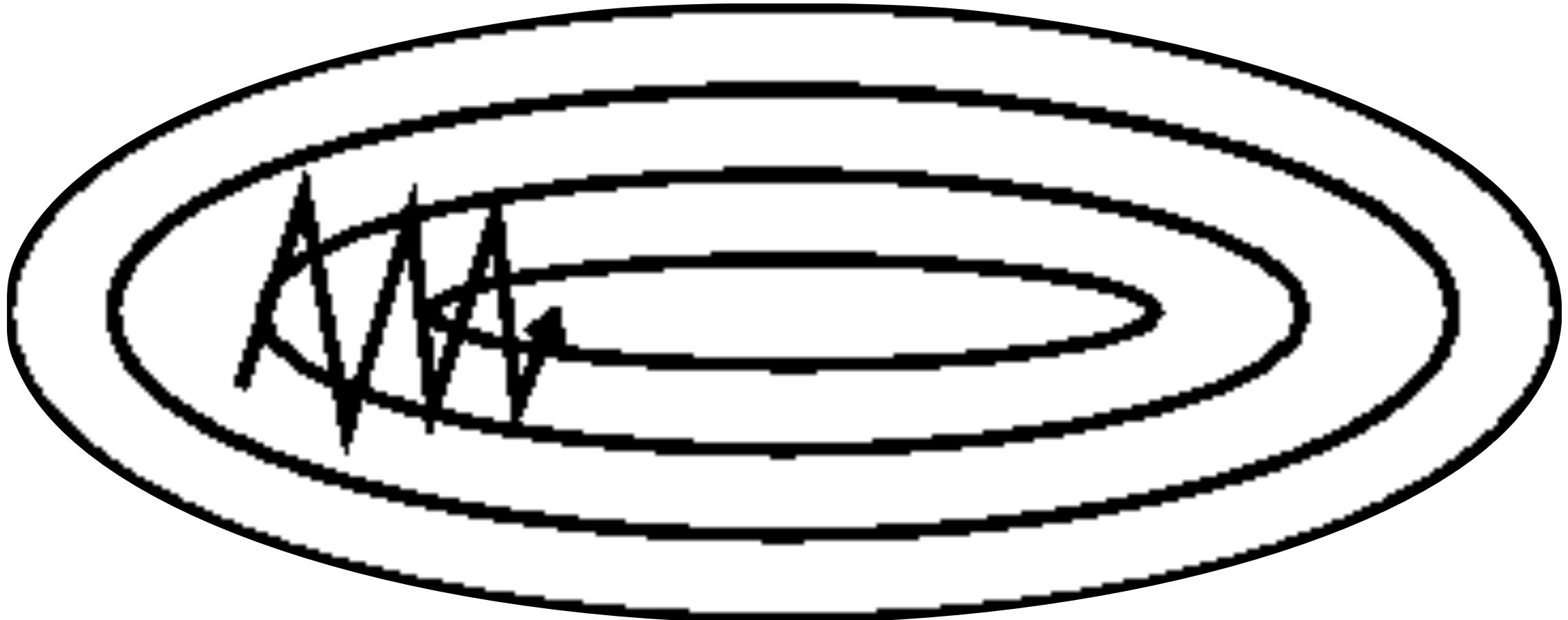
∇C = gradient of C

w = old weight

w' = new weight

η = learning rate (greek eta)

Gradient Descent and ravines



$$w \rightarrow w' = w - \eta \nabla C$$

∇C = gradient of C
 w = old weight

w' = new weight
 η = learning rate (greek eta)

Gradient Descent with Momentum

$$v \rightarrow v' = \mu v - \eta \nabla C$$

$$w \rightarrow w' = w + v'.$$

∇C = gradient of C

w = old weight

w' = new weight

η = learning rate (greek eta)

v = old velocity

v' = new velocity

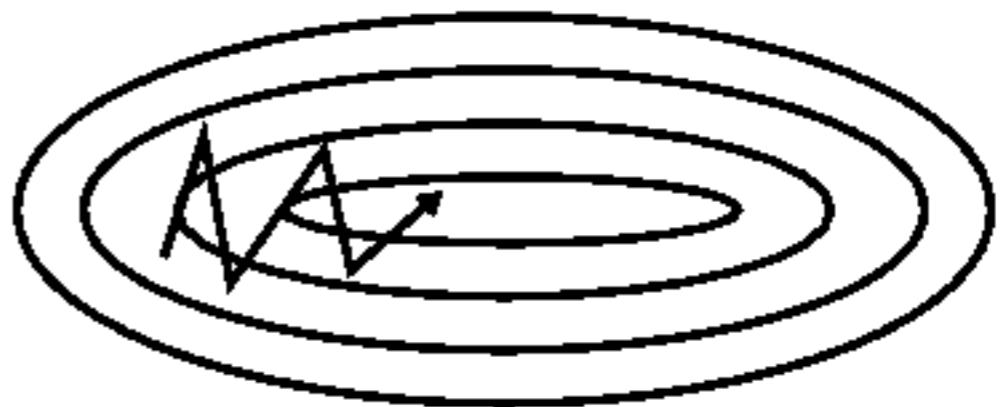
μ = momentum co-efficient or friction (greek mu)

Gradient Descent with Momentum

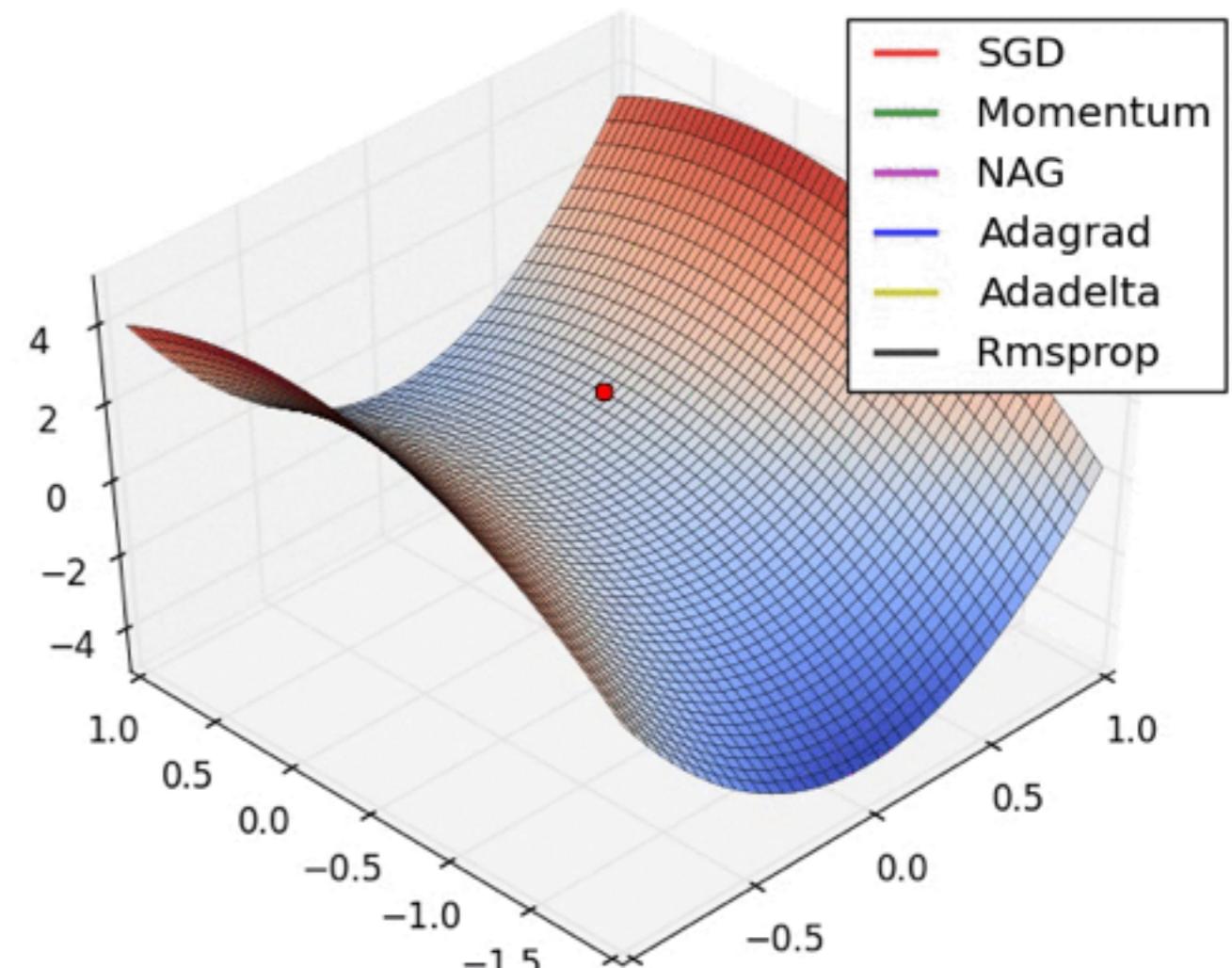
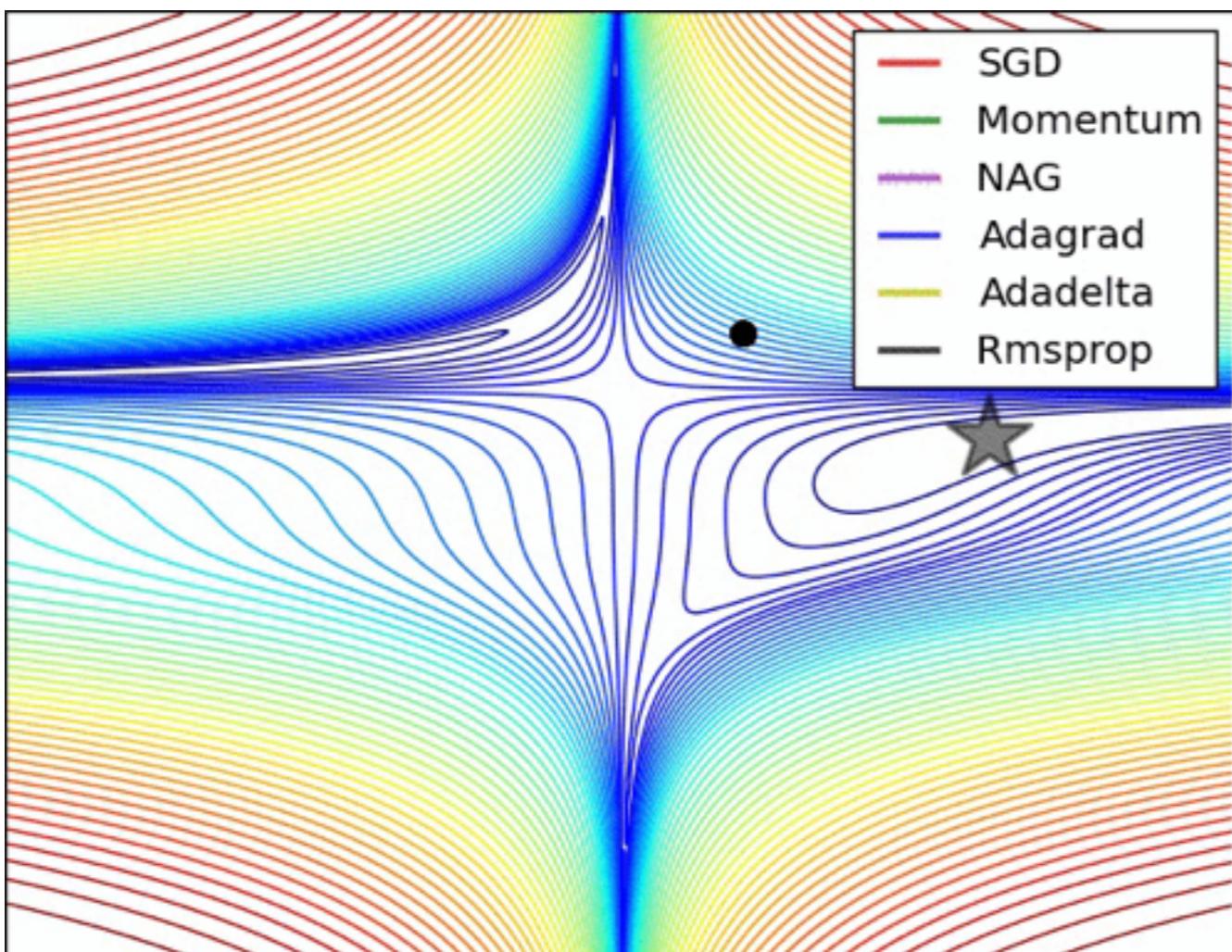
vanilla GD



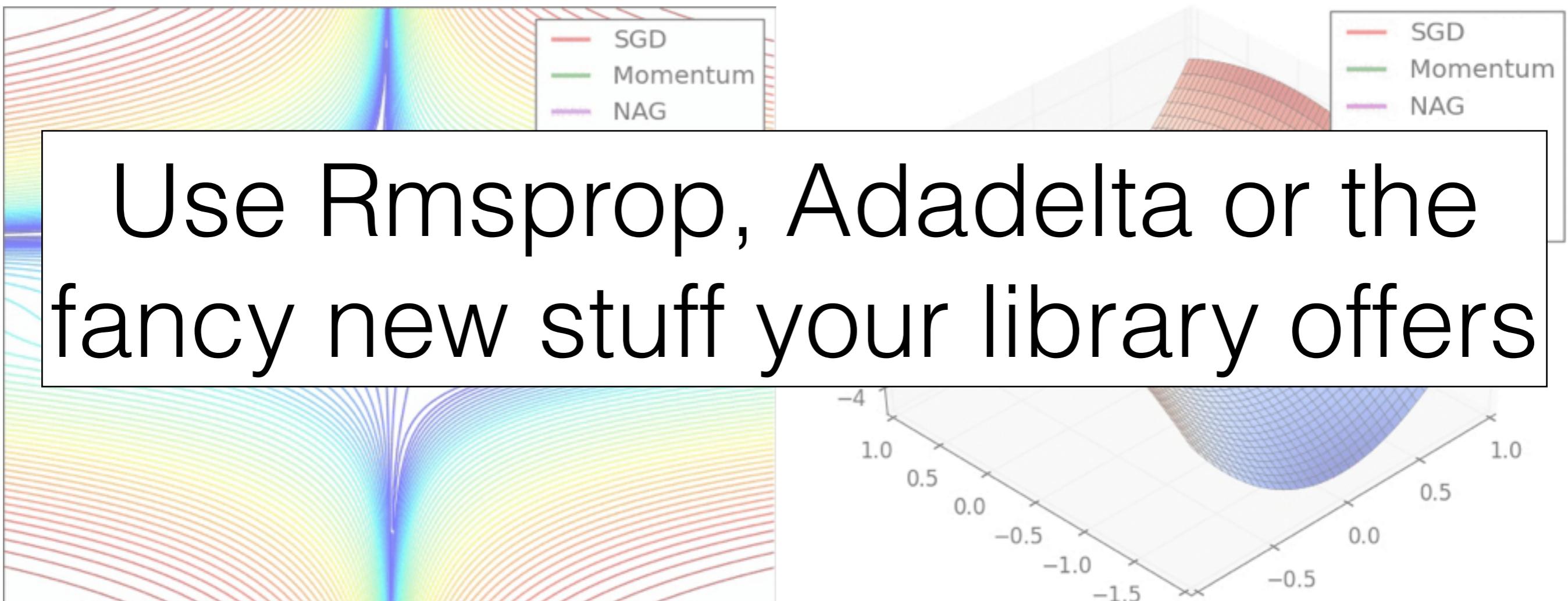
momentum SGD



But there are so many others...



But there are so many others...



back to

PROGRAMMING



Try out :

- dropout
- ReLUs
- **Different optimizers**
- **<https://keras.io/>**

`~/keras_mnist.ipynb`

Deep Learning and data sets

MNIST (1998)



- Hand written digits (0-9)
- Collected by the United States' National Institute of Standards and Technology, NIST
- Pre segmented
- 28x28 pixel in size
- 60k train, 10k test
- ('M' stands for modified)

MNIST 1998



- Hand written digits (0-9)
- Collected by the United States' National Institute of Standards and Technology, NIST
 - Follow
- Pre segmented
- 28x28 pixel in size
- 60k train, 10k test
- ('M' stands for modified)

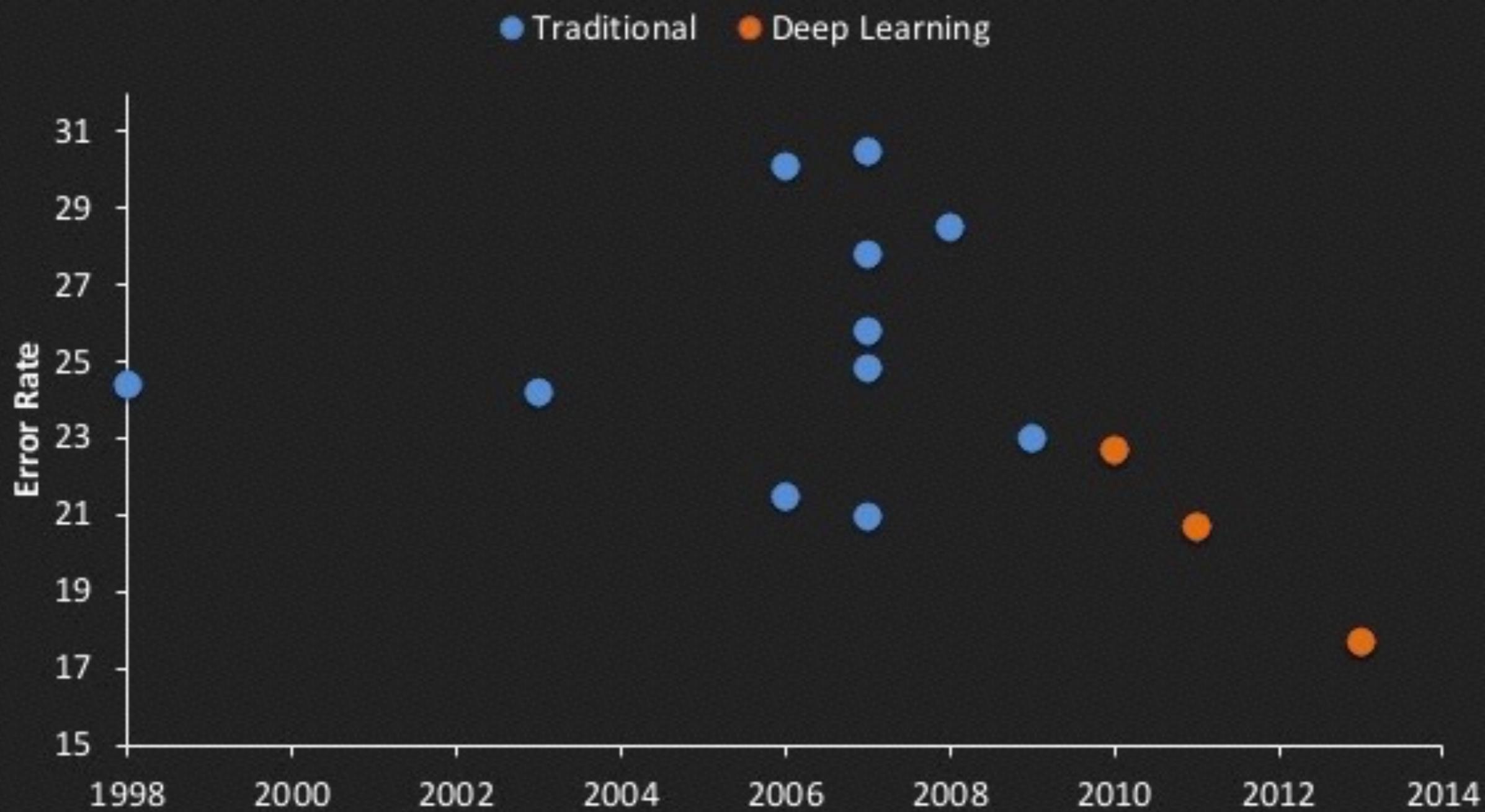
MNIST 1998

- Hand written digits (0-9)

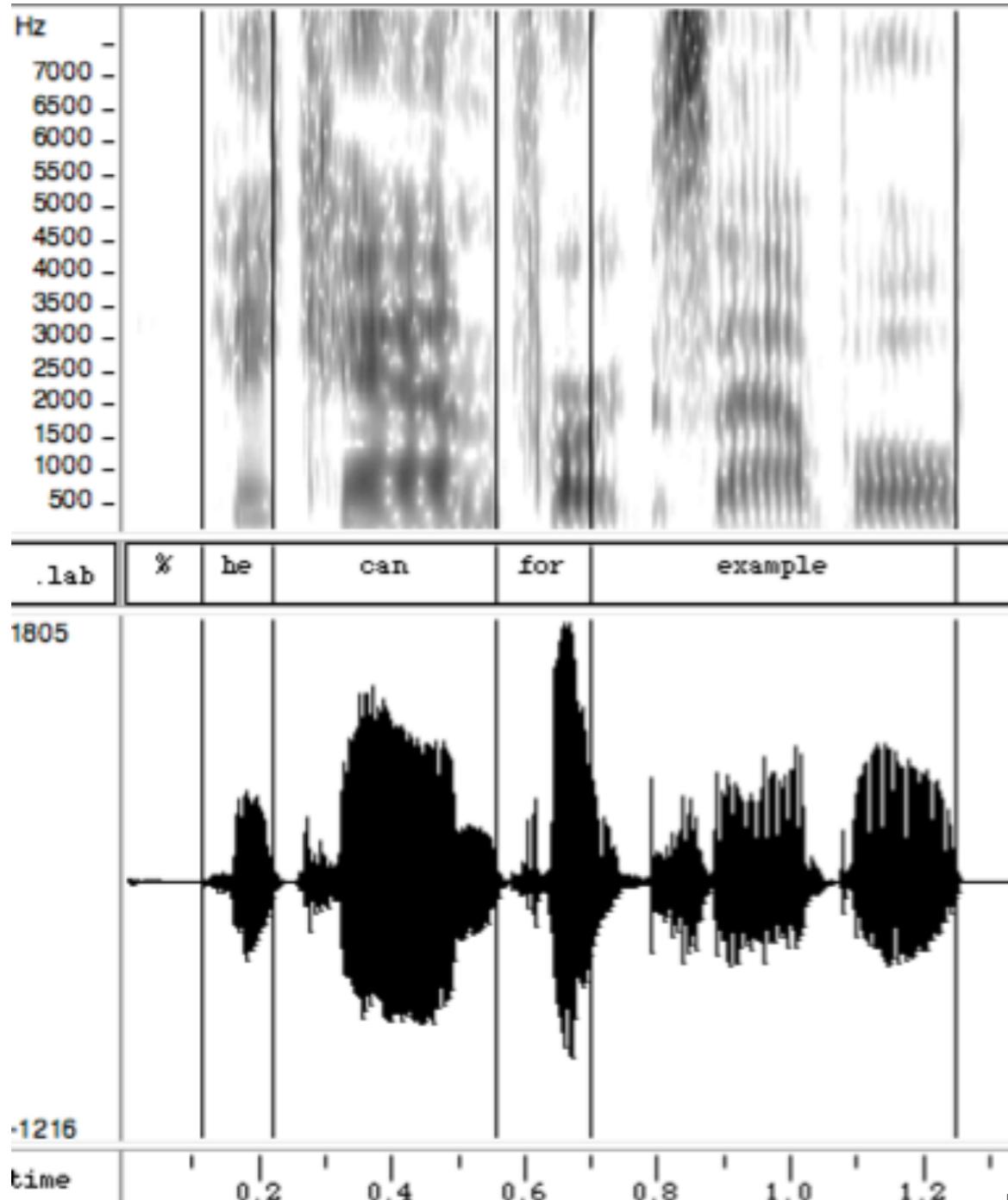


- 28x28 pixel in size
- 60k train, 10k test
- ('M' stands for modified)

TIMIT Speech Recognition

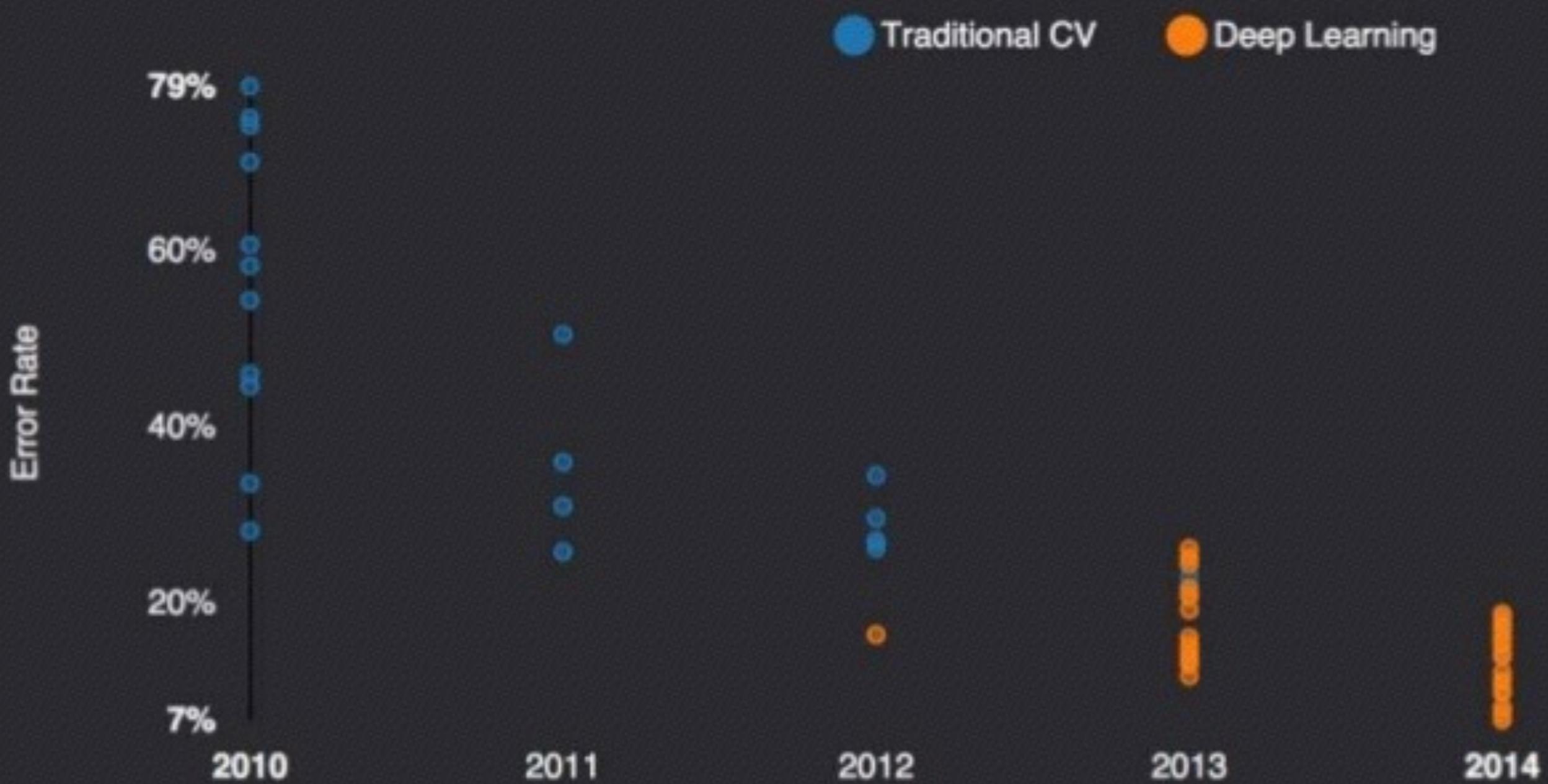


TIMIT (1990)

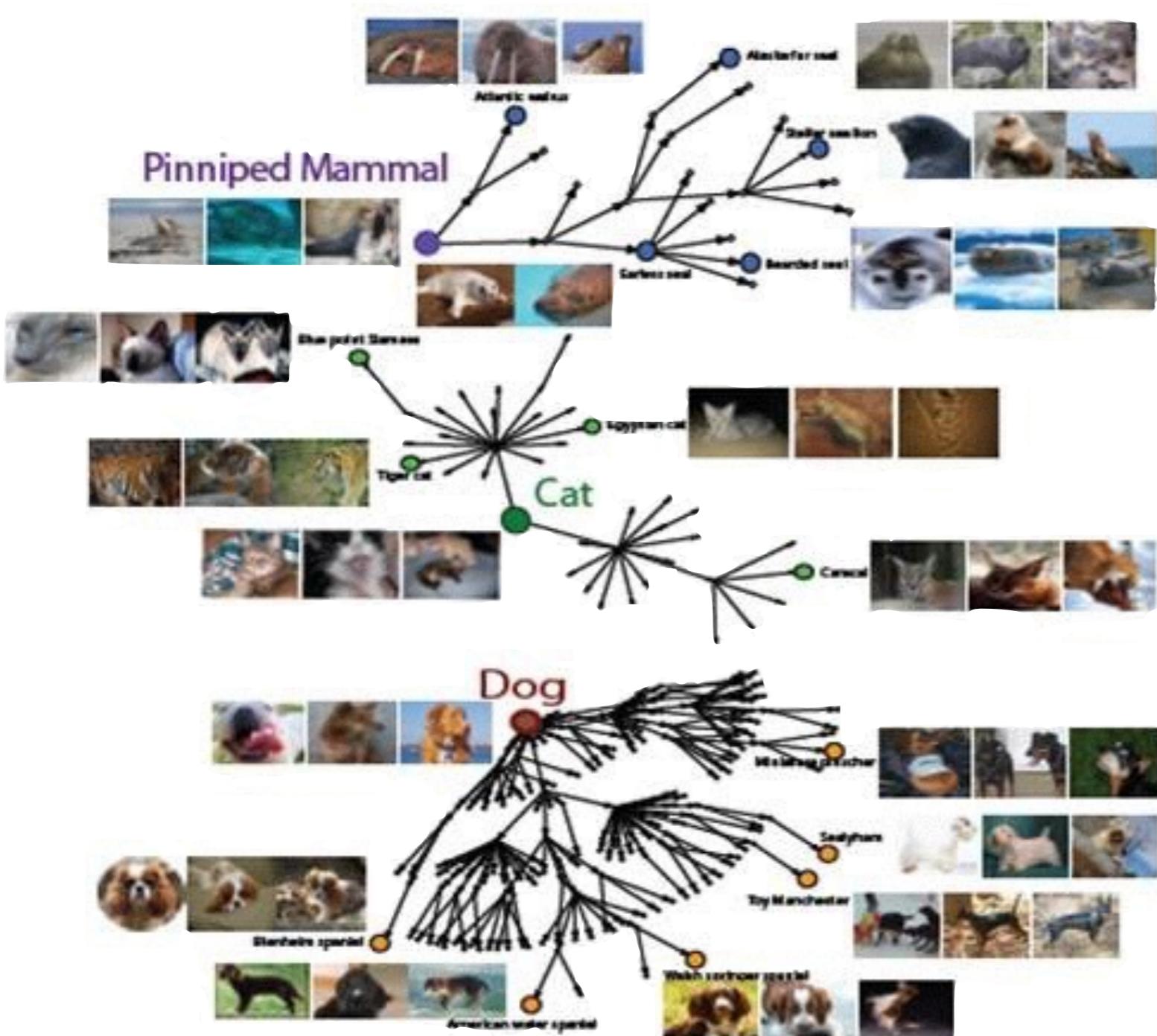


- Texas Instruments (TI) and Massachusetts Institute of Technology (MIT)
- 10 english sentences
- 630 speakers
- tagged on phone level
- 61 different phones
- 5.4h material
- @16 kHz with 16 bit
 - (~400m samples)
 - <https://catalog.ldc.upenn.edu/LDC93S1>

ImageNet Error Rate 2010-2014



'Imagenet' (2012)



Precise name: **Large Scale Visual Recognition Challenge**

- 1m images
- 1k concepts
- subset of the Imagenet
 - 20m images
 - 10k concept (tree structure)
- still growing
- image-net.org

'MS COCO' (2014)

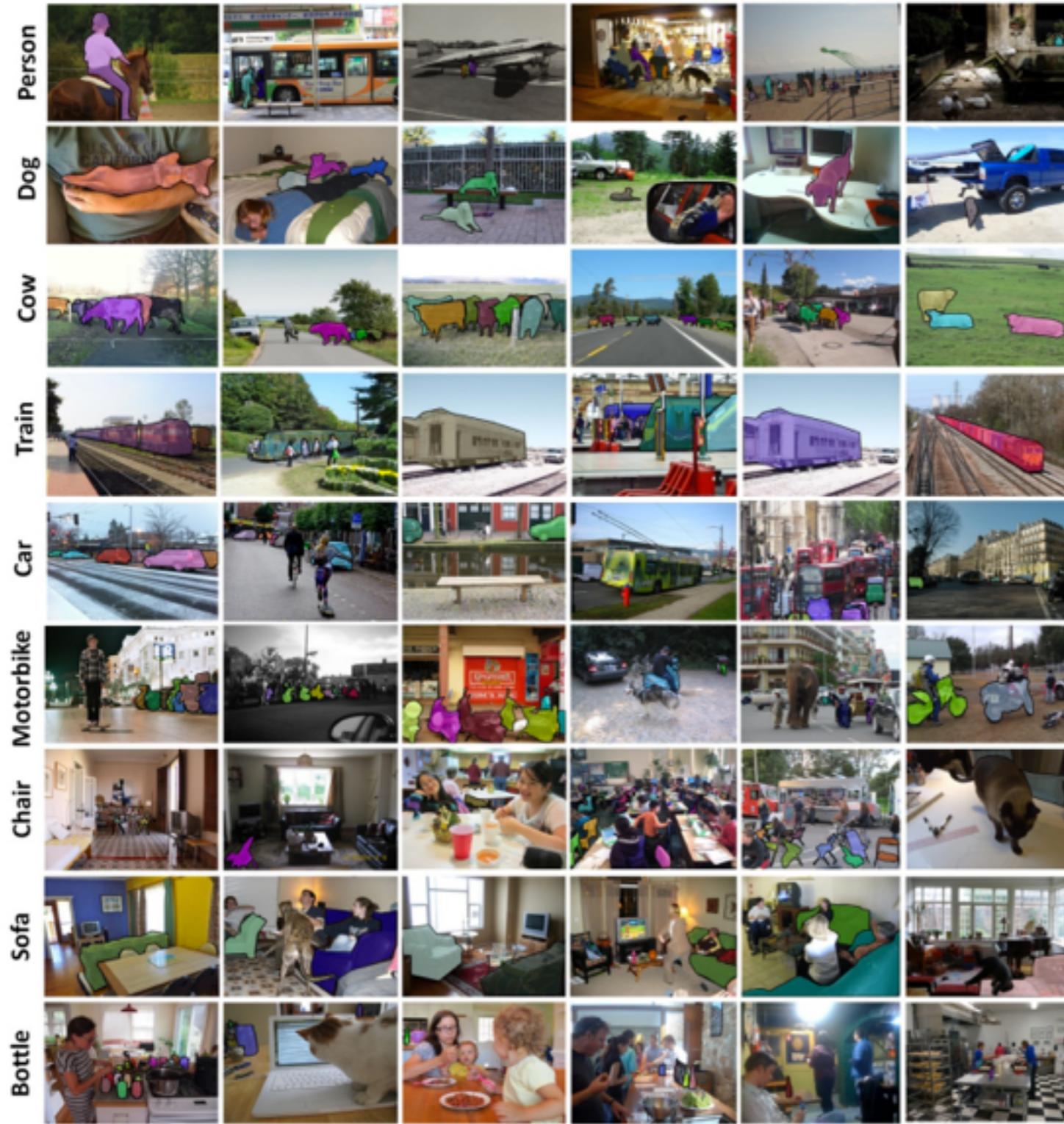
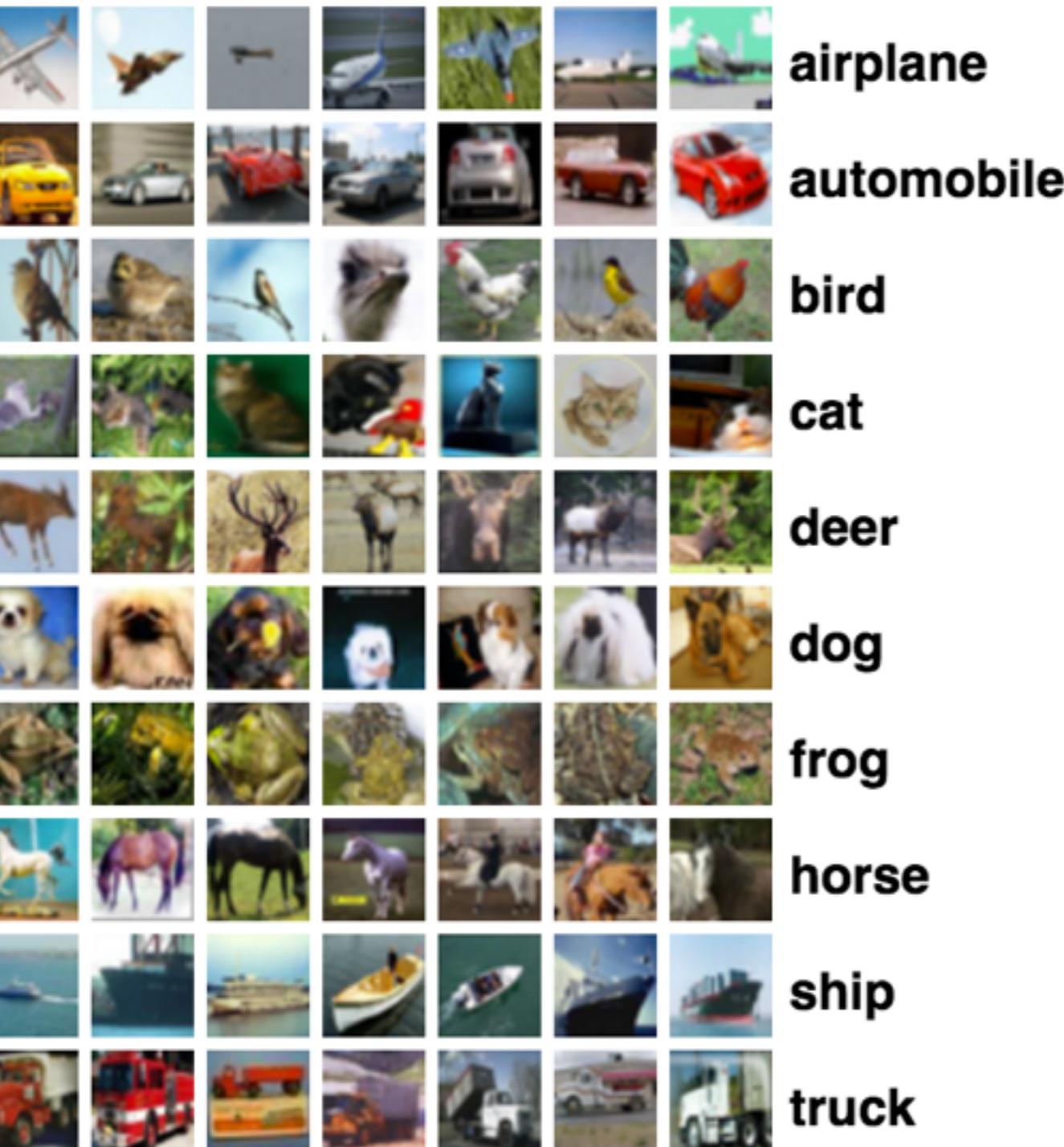


Fig. 6: Samples of annotated images in the MS COCO dataset.

Precise name: **Microsoft's Common Objects in COntext**

- More than 70 categories
- Object segmentation
- Multiple objects per image
- More than 300k images (150k)
- **5 captions per image**
- mscoco.org

CIFAR-10



airplane

automobile

bird

cat

deer

dog

frog

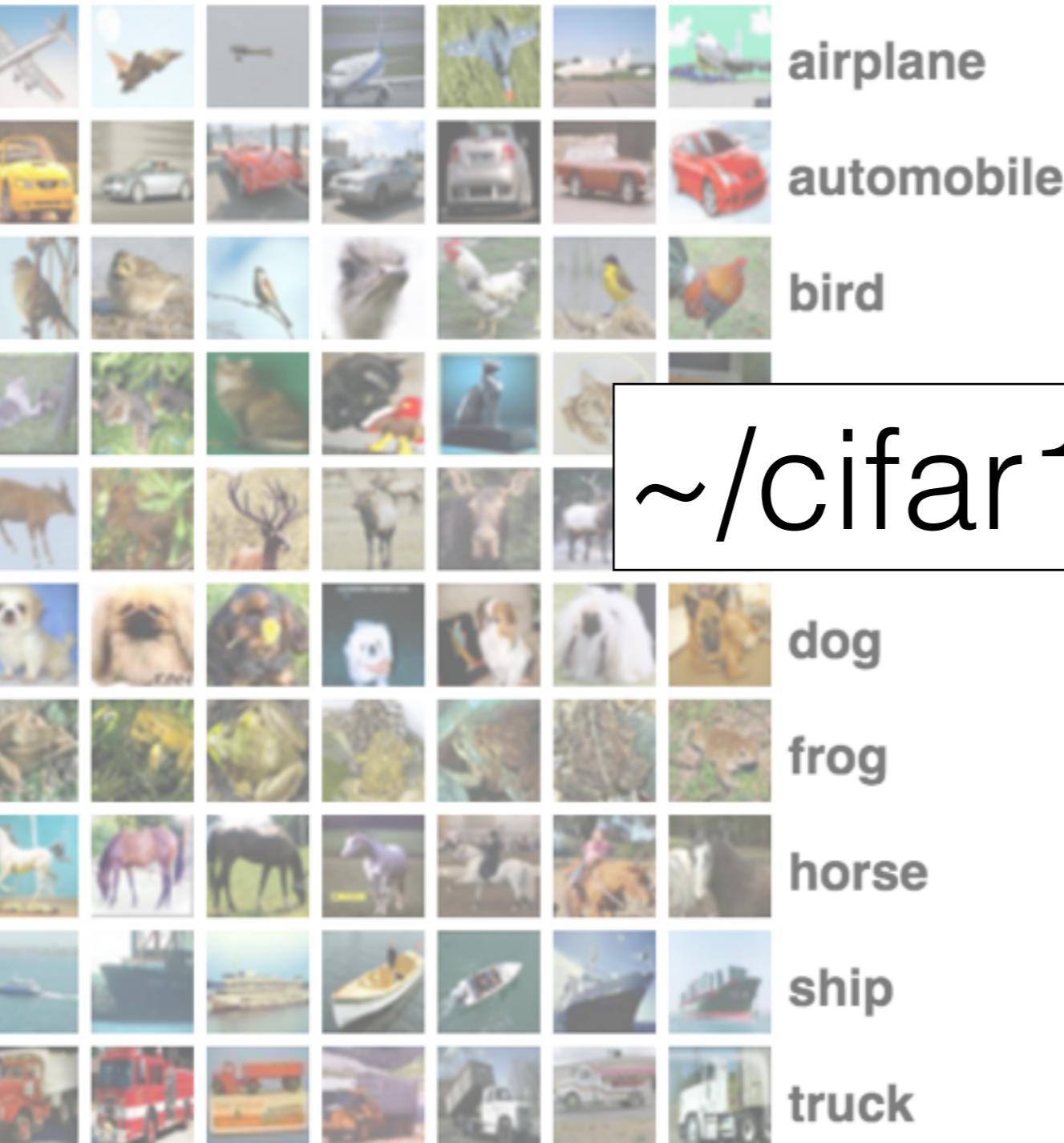
horse

ship

truck

- CIFAR stands for *Canadian Institute for Advanced Research*
- 10 concepts
- 6000 images per concept
- Image size is 32x32x3
- Part of the 80 Million Tiny Images dataset
- <https://www.cs.toronto.edu/~kriz/cifar.html>

CIFAR-10



- CIFAR stands for *Canadian Institute for Advanced Research*
 - 10 concepts
 - 6000 images per concept
- Part of the 80 Million Tiny Images dataset
- <https://www.cs.toronto.edu/~kriz/cifar.html>

~/cifar10.ipynb

2x32x3

Million Tiny Images

Cat or dog?

Or how to work with little data

Cat or dog?

Or how to work with little data

```
~/cats+dogs/train/cats/...
```

from: www.kaggle.com/c/dogs-vs-cats

Cat or dog?

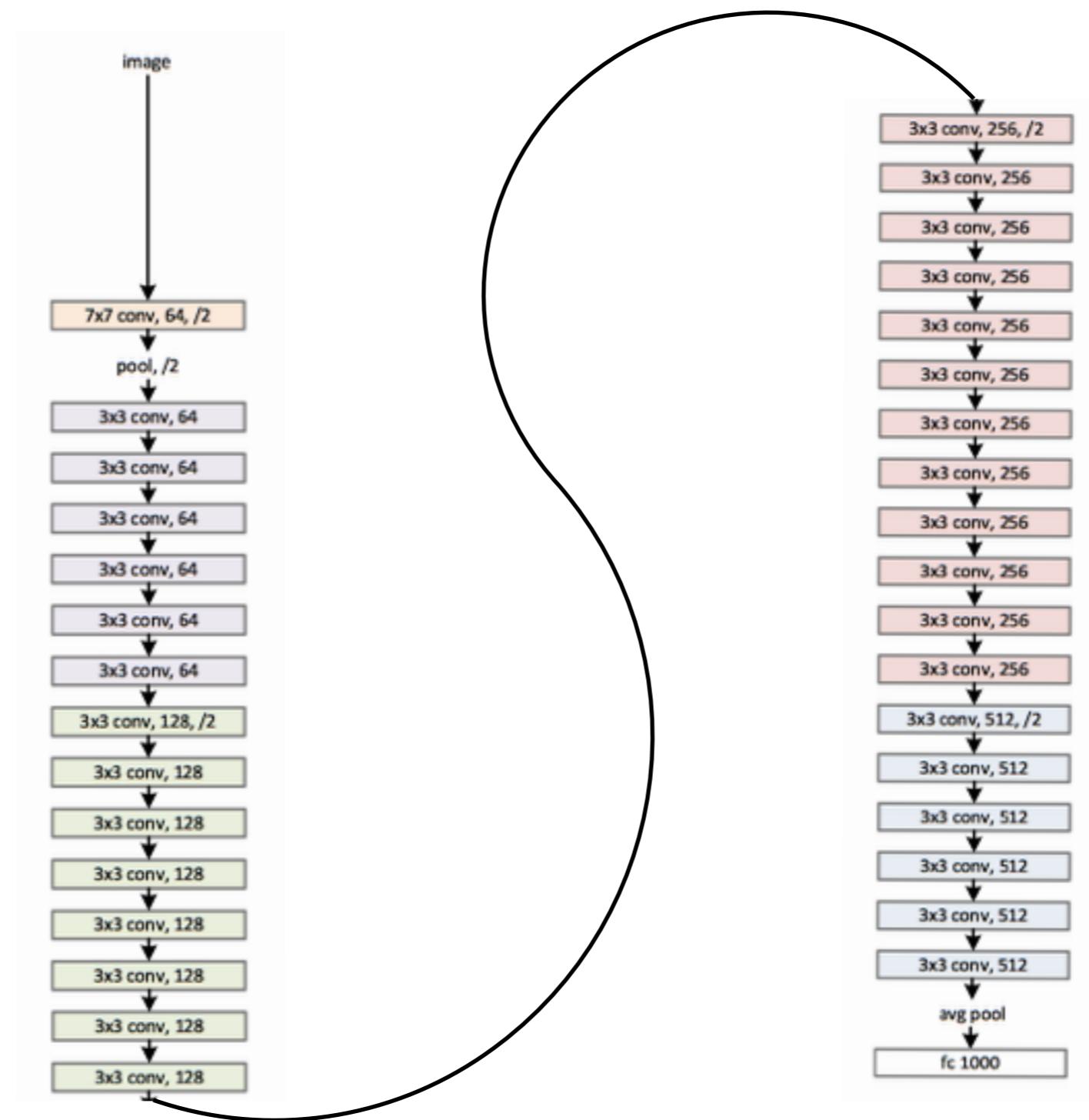
Or how to work with little data

```
~/cats+dogs/first_try.ipynb
```

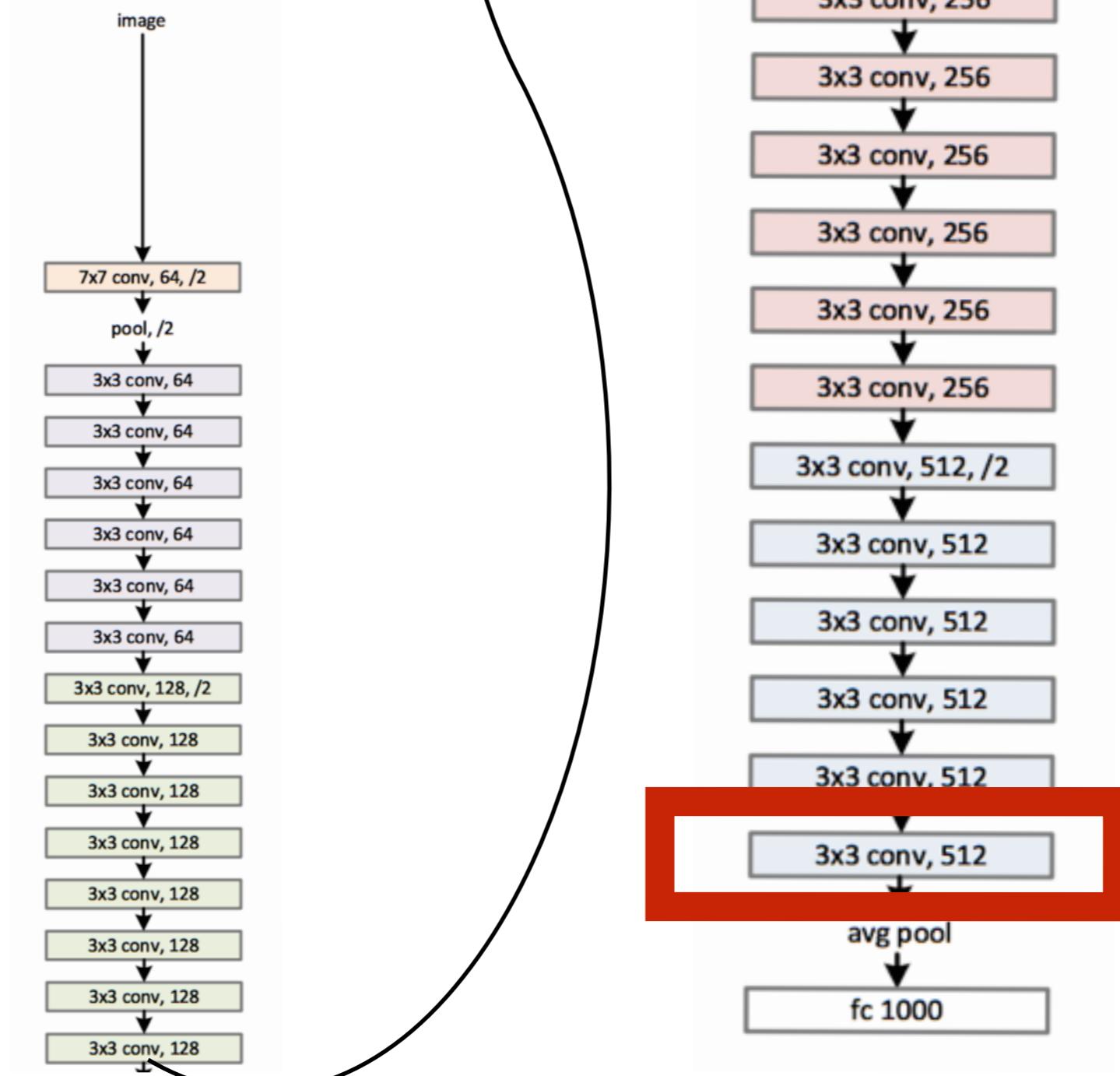
from: www.kaggle.com/c/dogs-vs-cats

The trick

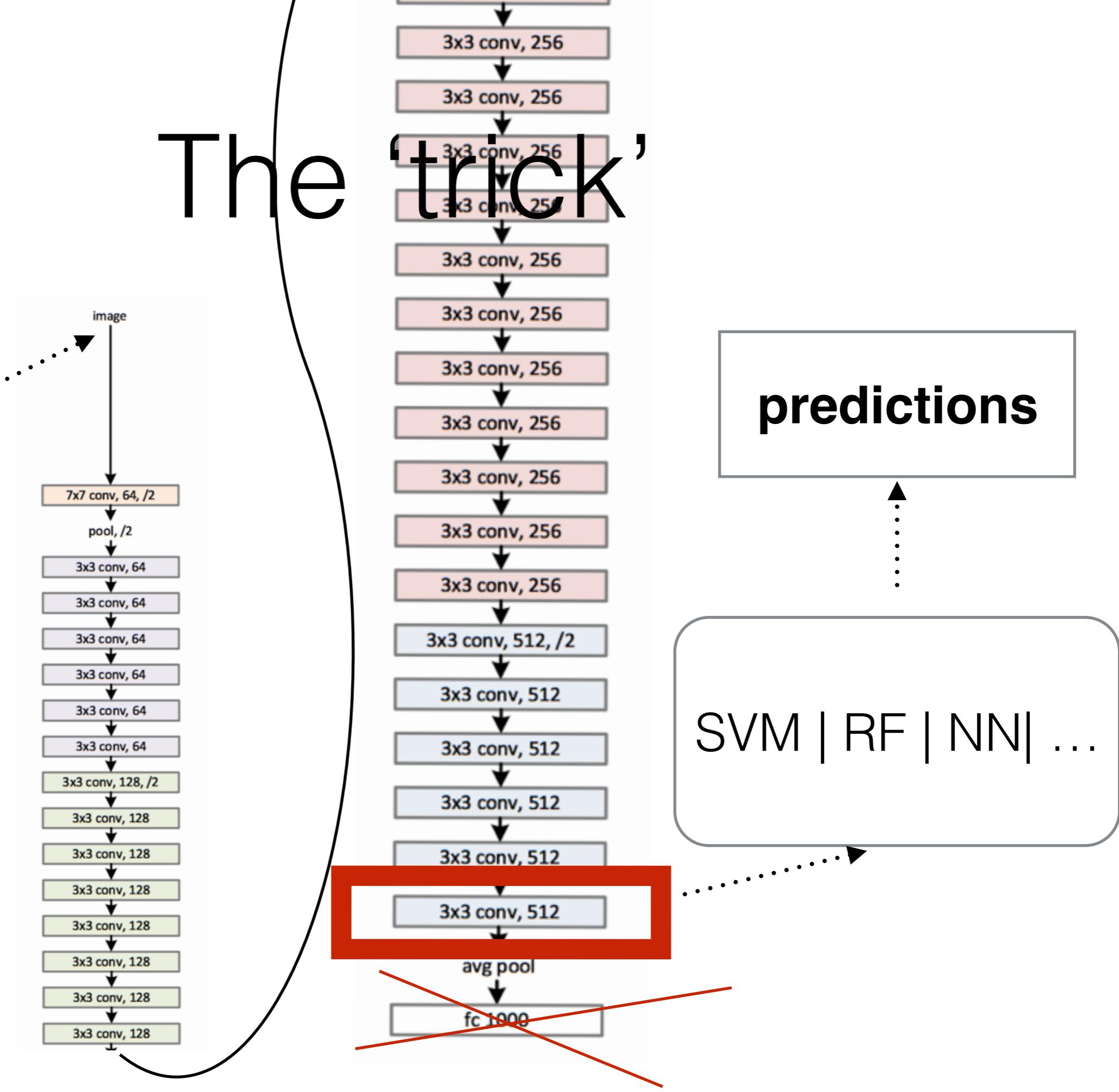
The ‘trick’



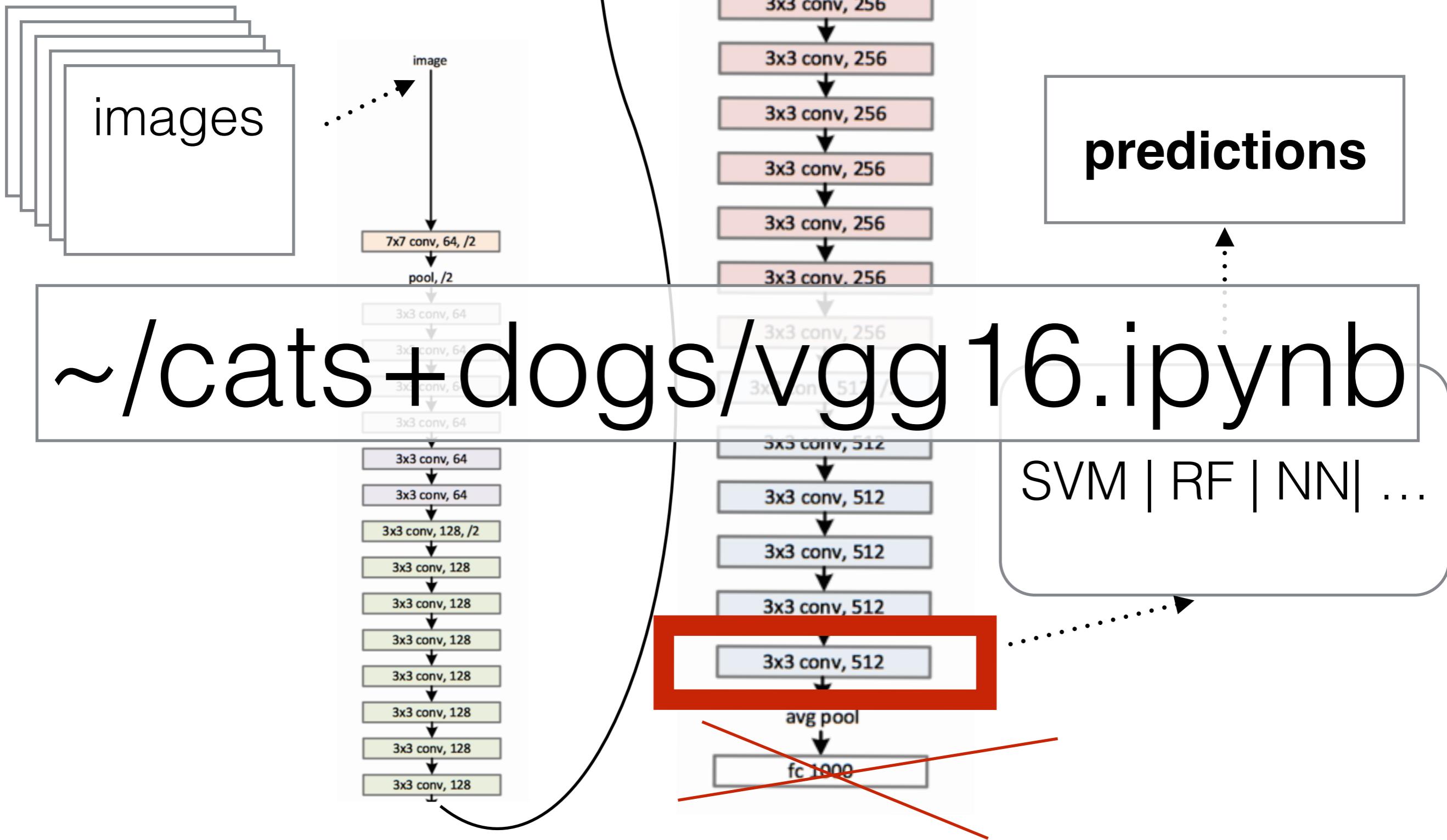
The 'trick'



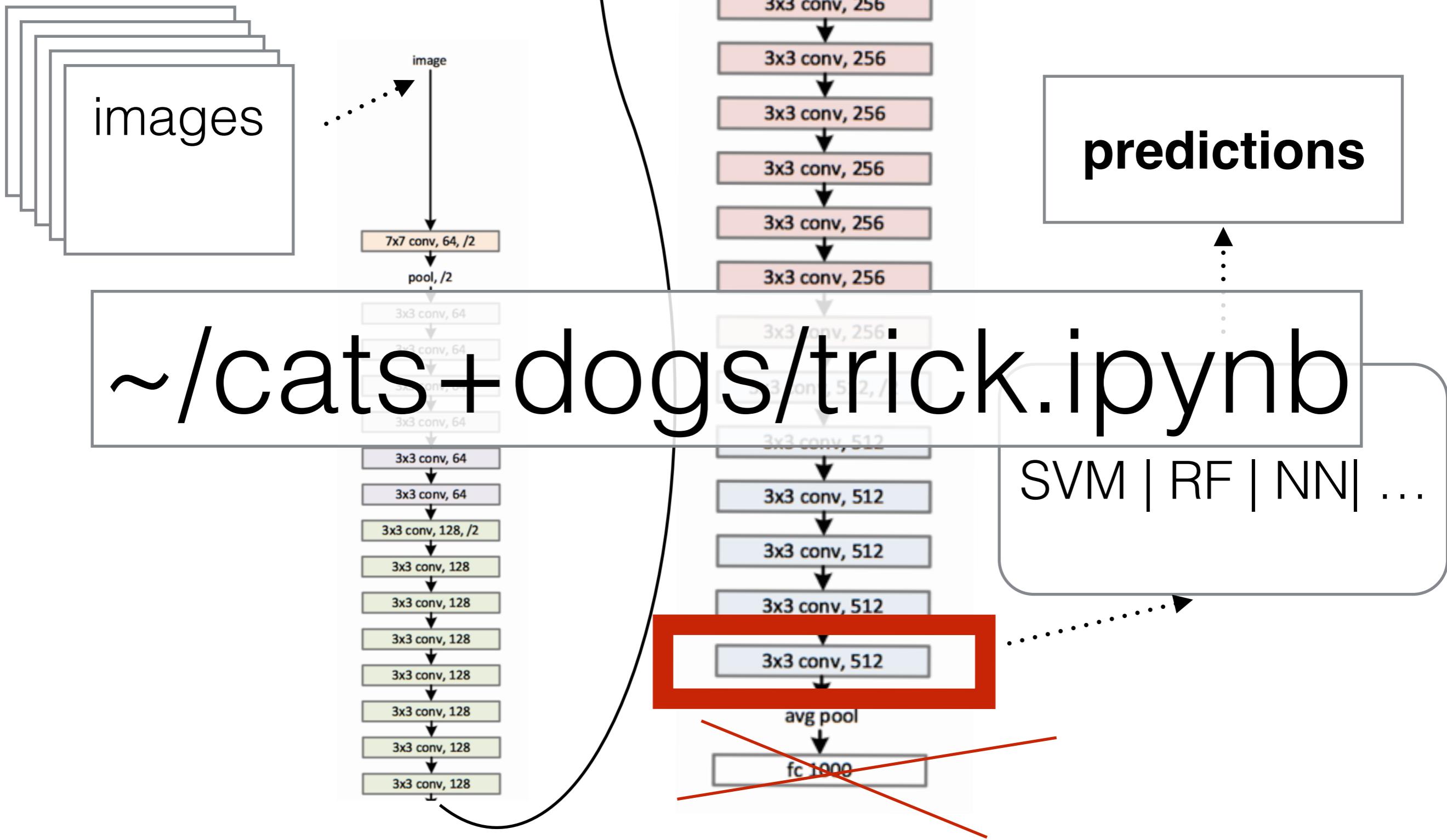
images



The 'trick'



The 'trick'



The ‘trick’

Using pre-trained weights of a feature extractor in combination with a familiar classifier like SVM or RF.

Pros

- predictable training time
- learned features
- needs less data

Cons

- needs pre-trained network weights
- introduction of bias
- reaches not the full potential of Deep Learning

The ‘trick’

Using pre-trained weights of as feature extractor in combination with an familiar classifier like SVM or RF.

Pros

`~/cats+dogs/trickier.ipynb`

- learned features
- needs less data

Cons

`~/cats+dogs/trickier.ipynb`

- predictable training time
- needs pre-trained network weights
- introduction of bias
- reaches not the full potential of Deep Learning

Q: How do I know what architecture to use?

Q: How do I know what architecture to use?

A: Don't be a hero:

1. Take whatever works on ILSVRC (latest ResNet)
2. Download a pretrained model
3. Potentially add/delete some parts of it
4. Finetune to your application

Q: How do I know what architecture to use?

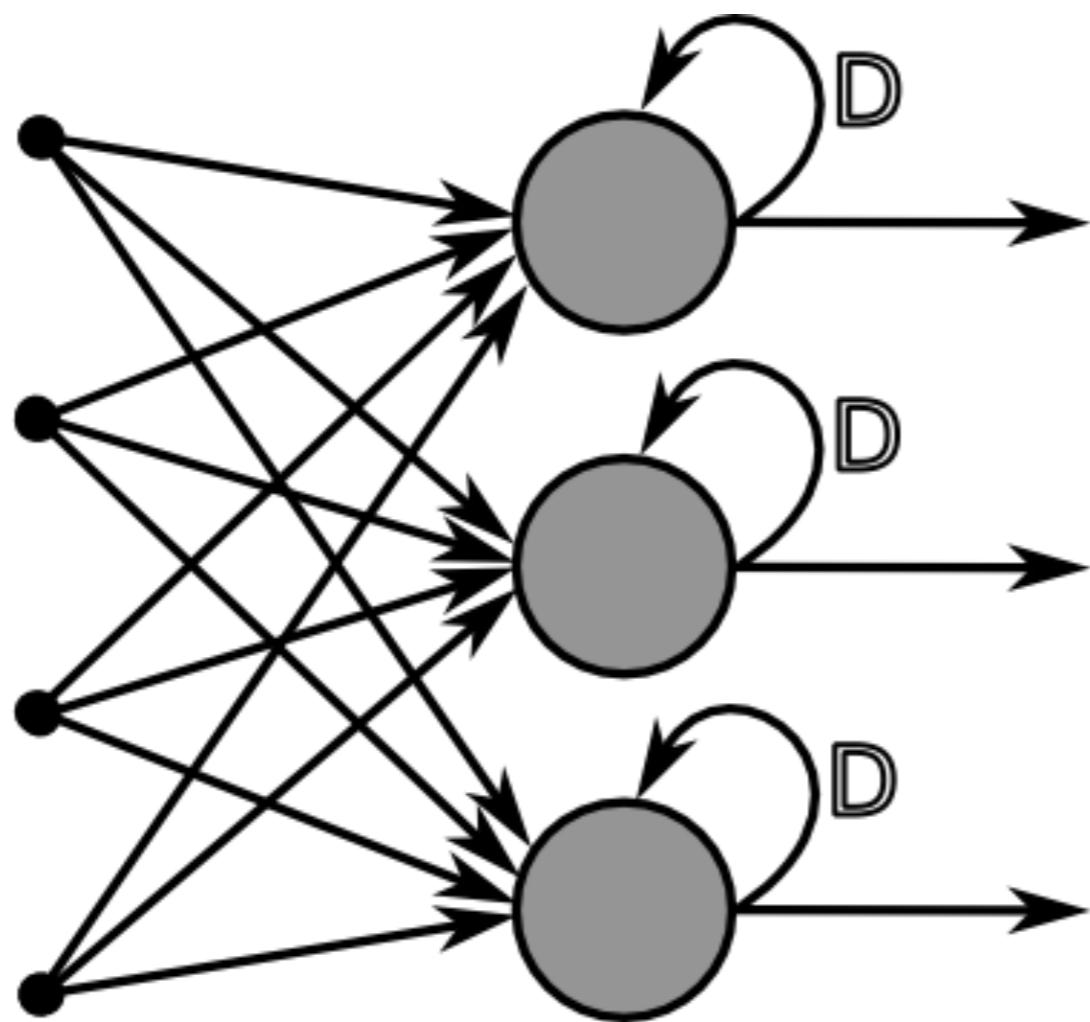
A: Don't be a hero:

1. Take whatever works on ILSVRC (latest ResNet)
2. Download a pretrained model
3. Potentially add/delete some parts of it
4. Finetune to your application

***Architecture engineering
is the new Feature engineering***

Other Deep Learning approaches

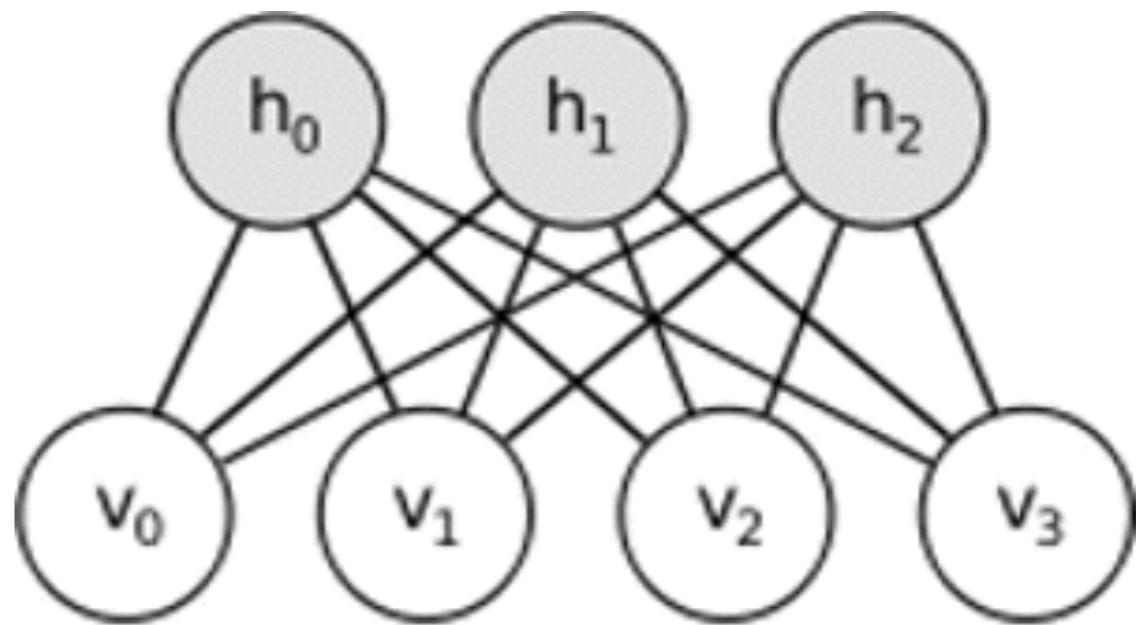
Recurrent Neural Networks (RNN)



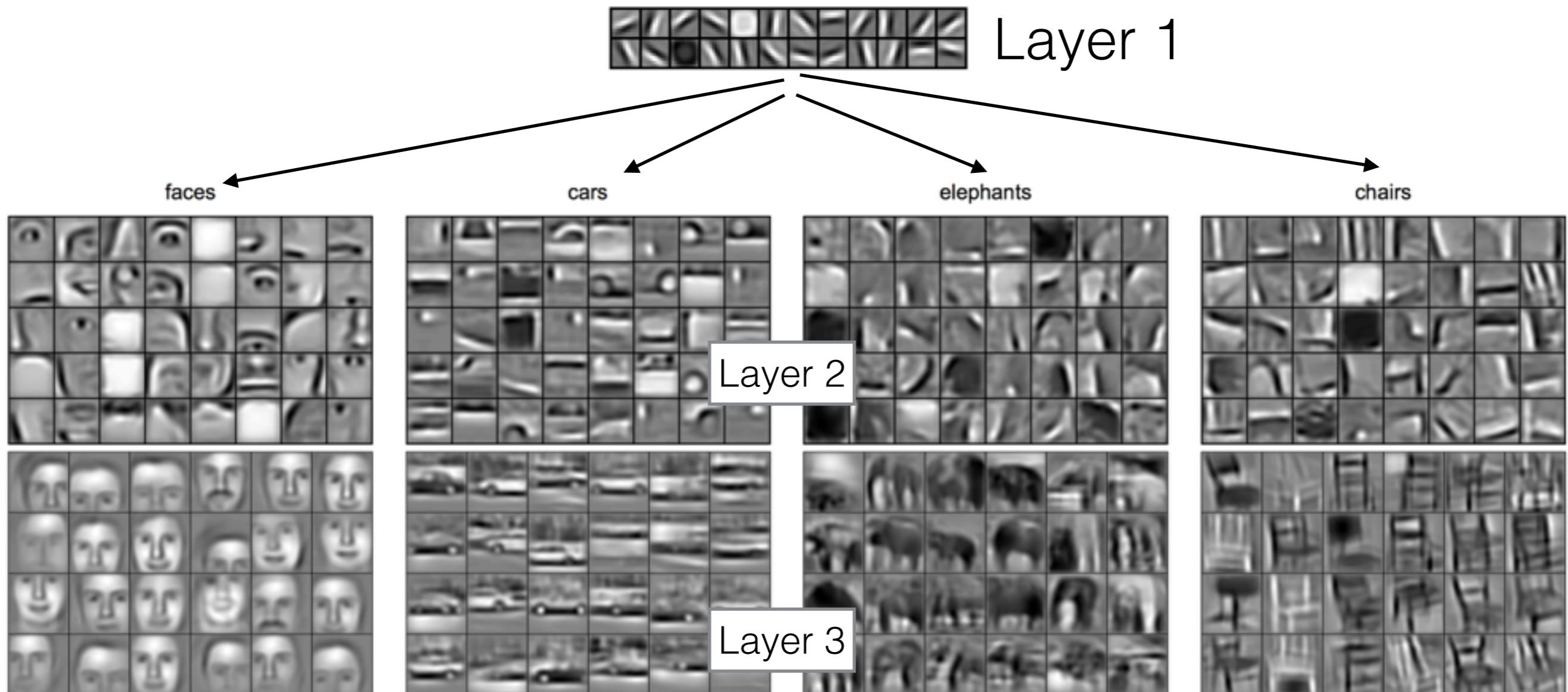
- Good for time series problems
 - sound
 - hand writing
- unstable
 - vanishing and exploding gradients
- but recent advancement
 - with major breakthroughs in voice recognition (google now, siri, cortana)

Restricted Boltzmann Machines (RBMs)

- are also known as **Deep Believe Networks**
- train a reconstruction of the inputs with fewer parameters
- can be trained unsupervised or semi-supervised
- were the first general-purpose deep architecture that could be trained successfully (2006)
 - give easily interpretable nodes
 - are now surpassed by CNNs



Convolutional Deep Believe Networks



Paper: [Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations](#)

Deep Learning small talk

The Deep Learning Heads



Yan LeCun



- Known for CNNs
- LeNet5 1989 (AT&T)
- Torch / Lua
- now head of Facebook AI Research (working there also with V. Vapnik (SVM))

Geoffrey Hinton



- Co-inventor of the backpropagation in '74
- Known for DBNs
- University Toronto + 50% google research since 2012
- Vision break through 2012

Yoshua Bengio



Yoshua Bengio

- University Montreal
- Sep 2016: Grant of \$93 million to work on "AI"
- Theano

Jürgen Schmidhuber



- Pronounciation:
You_again Shmidhoobuh
- Known for RNN
- University of Lugano
(Switzerland)
- Also working on art theory
+ generation through NNs

The Neural Computation and Adaptive Perception program



- founded by LeCun, Hinton and Bengio in 2004
- funded by a Canadian Institute for Advanced Research (CIFAR) grant
- invite only
- laid the groundwork and spearheaded the Neural Network renaissance

The Neural Computation and Adaptive Perception program



Deep Learning Conspiracy



VS



OpenAI

- non-profit artificial intelligence research company
- founded by 2015 by Elon Musk and Sam Altman (YCombinator)
- ‘1 Billion \$ in funding’

“Carefully promote and develop friendly AI in such a way as to benefit, rather than harm, humanity as a whole.”

OpenAI

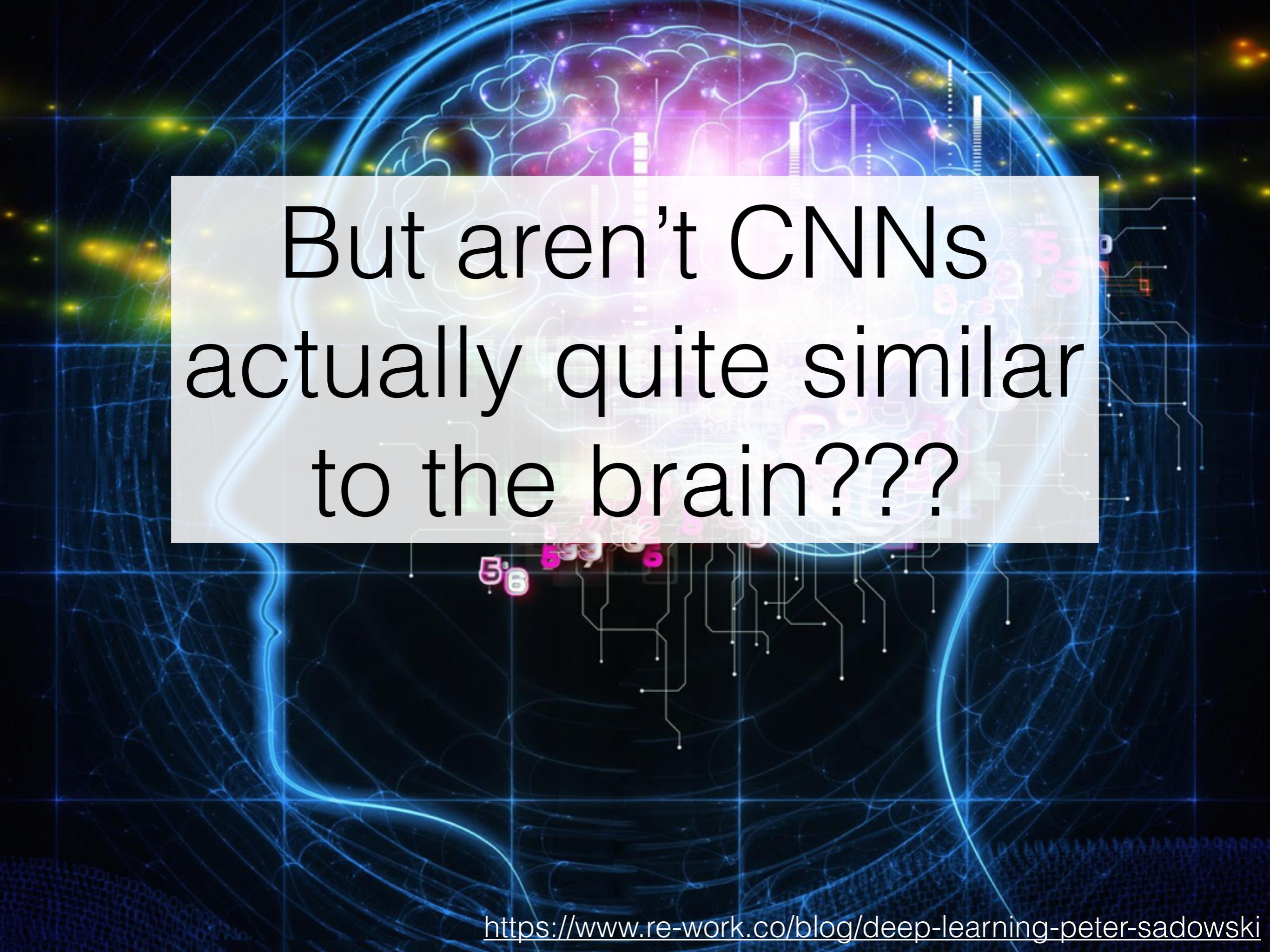
- non-profit artificial intelligence research company

OpenAI's Director of 'Strategy and Communications', Jack Clark, has a very good newsletter about Deep Learning and AI.

Follow it to stay up to date:

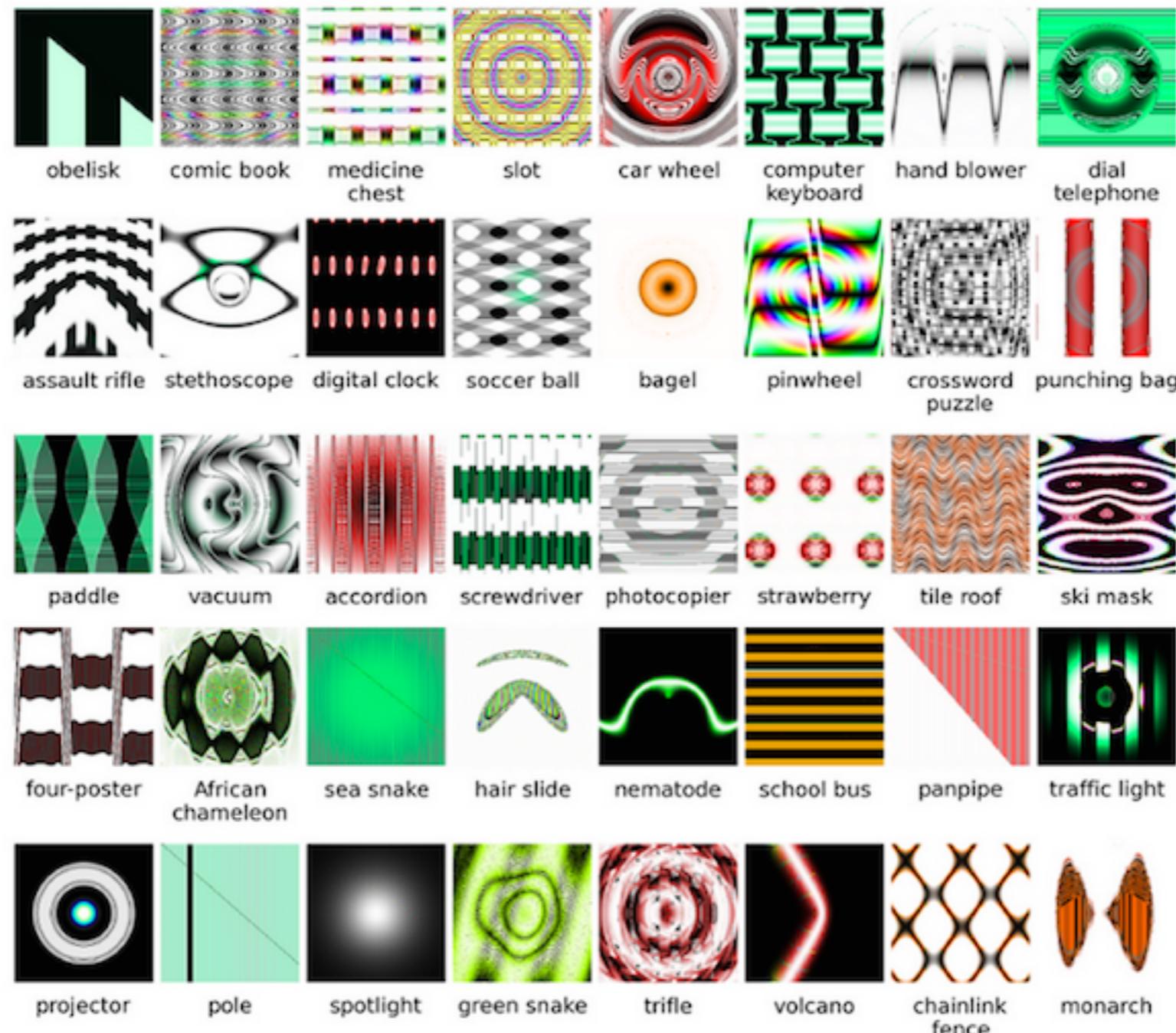
jack-clark.net

wnoie.

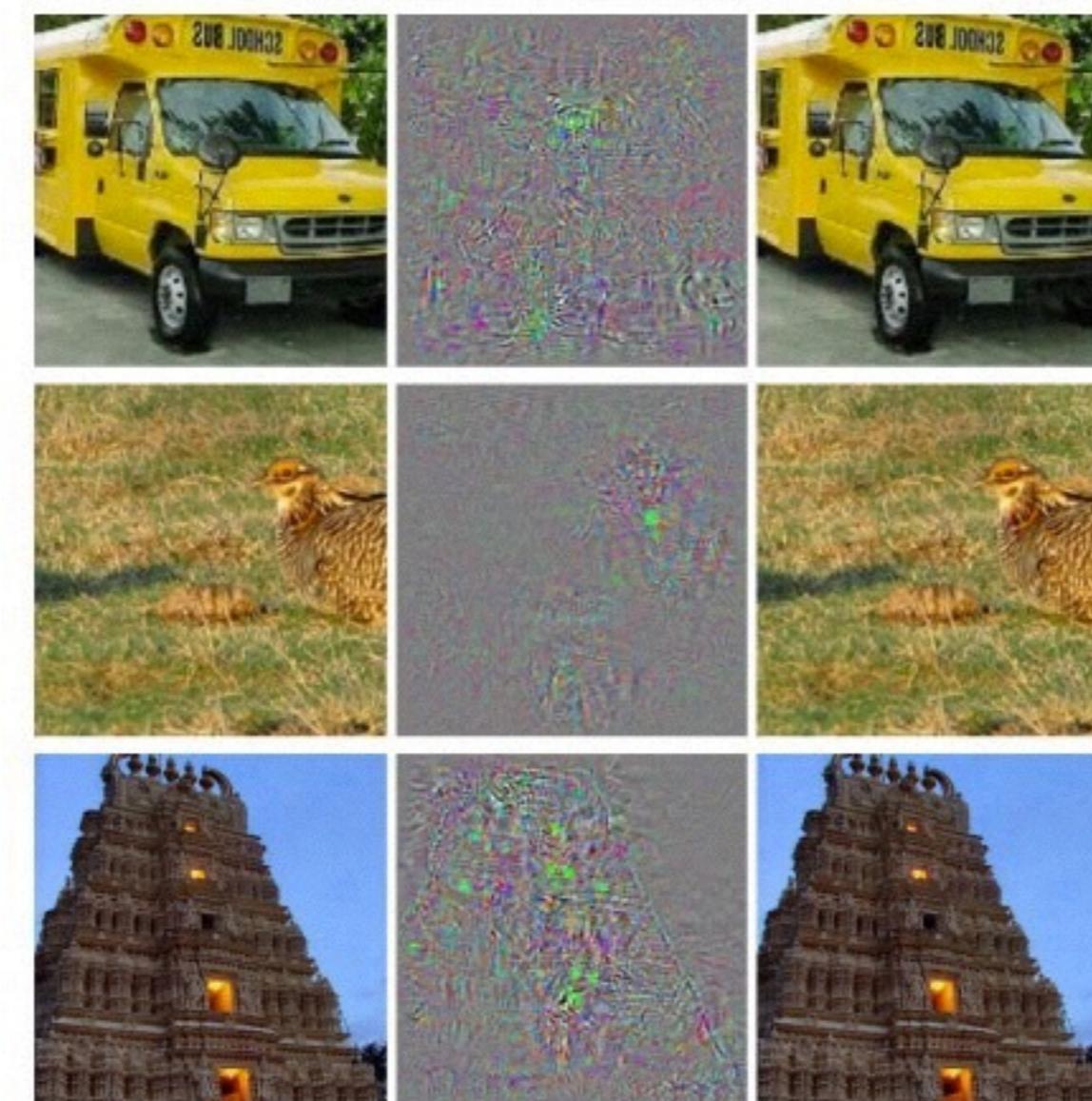


But aren't CNNs
actually quite similar
to the brain???

Adversarial Examples



Adversarial Examples

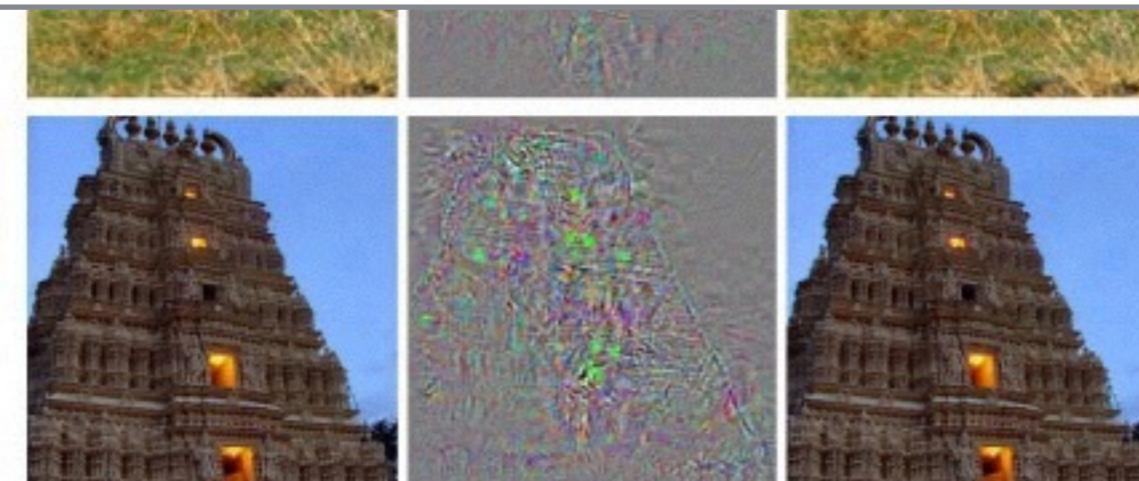


Paper: Intriguing properties of neural networks

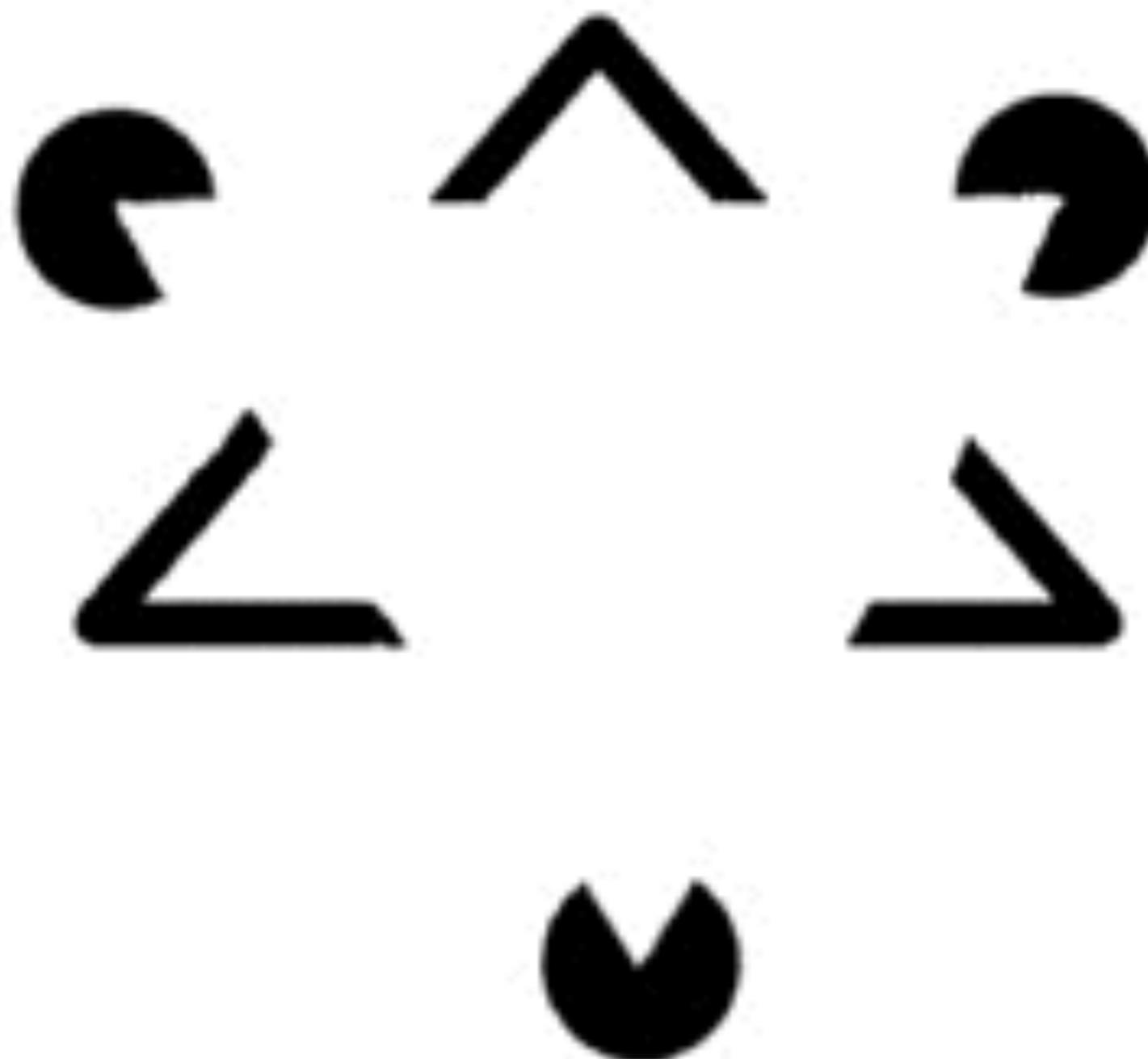
Adversarial Examples



More training data?!



Gestalt Law of Closure



Conclusion

Pros

- no feature engineering necessary
- setting state of the art for
 - computer vision,
 - speech recognition and
 - some text analysis
- full potential still unknown
 - they get better with more data

Cons

- doesn't work out of the box
- many choices
 - architecture
 - hyper parameter
 - learning algorithm
- needs huge amounts of data
- computational expensive
- lots of expensive experiments

Thank you